

1

2

3

4

5

6

7

SUPPLEMENTARY INFORMATION TO THE MANUSCRIPT

8

CHEMICAL UNCLONABLE FUNCTIONS BASED ON OPERABLE RANDOM DNA POOLS

9

10 Anne M. Luescher¹, Andreas L. Gimpel¹, Wendelin J. Stark¹, Reinhard Heckel² und Robert N. Grass^{1*}

11 ¹ Department of Chemistry and Applied Biosciences, ETH Zürich, Vladimir-Prelog-Weg 1-5, 8093,
12 Zürich, Switzerland

13 ² Department of Computer Engineering, Technical University of Munich, Arcistrasse 21, 80333,
14 Munich, Germany

15 *Corresponding author, robert.grass@chem.ethz.ch

16

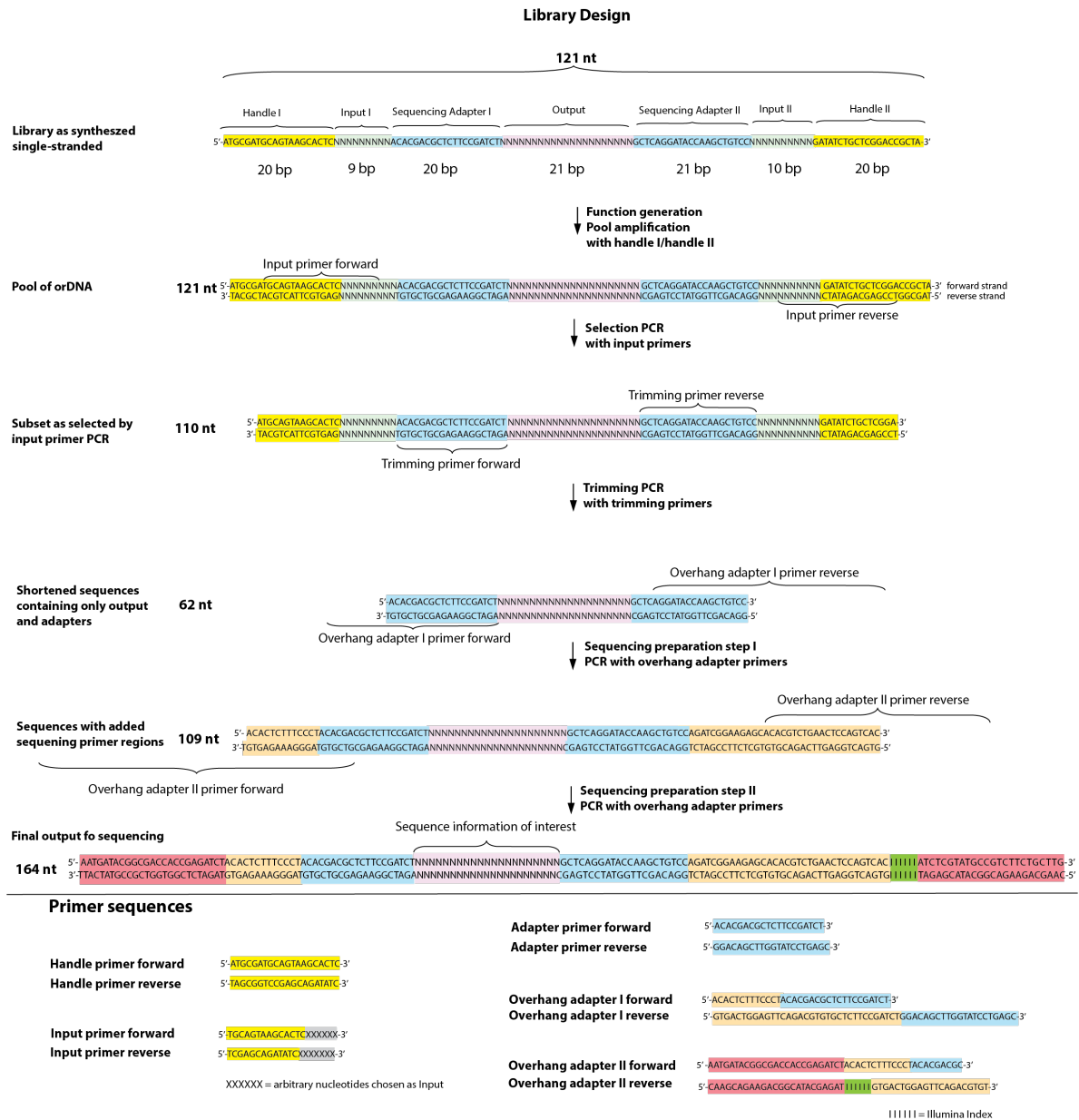
17 **Supplementary Information**

18 Supplementary Figures 1-17

19 Supplementary Notes 1-13

20 Supplementary Tables 1-5

21



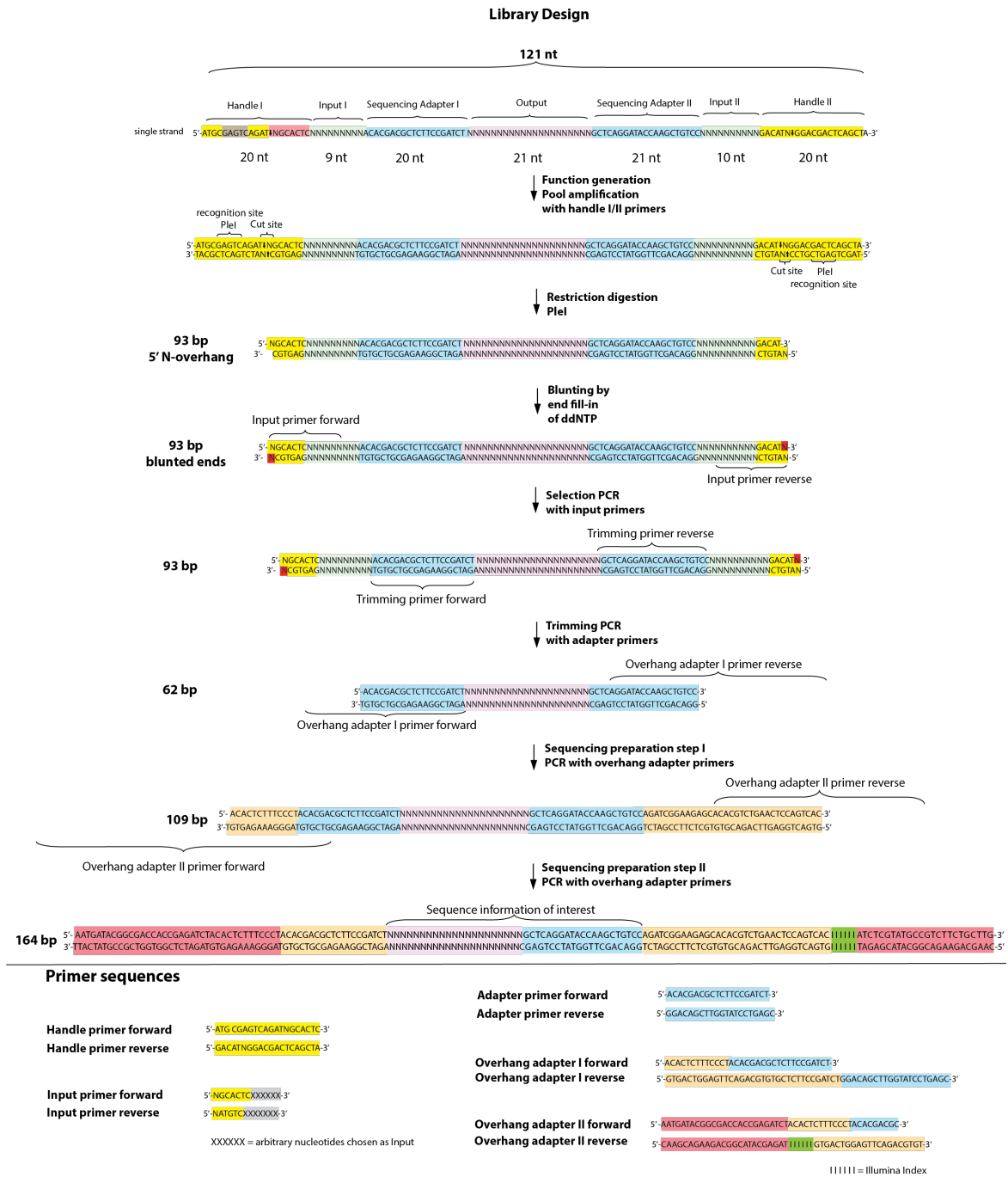
22

23 **Supplementary Fig. 1: Library design**

24 Shown are the sequence structures at various stages of the process of creating and operating the
 25 chemical function. First, a pool comprising several milligrams of single-stranded
 26 oligonucleotides is obtained from solid-state synthesis, comprising both sequence-determined and
 27 random segments. A subset of this library (e.g. $10^8 - 10^{11}$ sequences) is then amplified via PCR using
 28 the outer handles as primers. The orDNA pool can then be operated using a set of input primers. The
 29 two inputs comprise a few bases (number may vary between 6 and 9), which selectively bind to the
 30 sequences in the pool with complementary matching segments in their randomly synthesized regions.
 31 Aside from the bases binding to the input regions, the primer further contains a part of the handle
 32 sequences to guide the input towards the correct binding position and to add enough length to the
 2

33 primers for successful amplification. The selected sequences are only a very small subset of the
34 orDNA pool and are identified via next generation sequencing. To prepare sequencing, three further
35 PCR reactions are conducted. First, the sequences are trimmed using the two adapter primers. The
36 vast majority of the resulting pool then only contains the two adapters with the output section in
37 between. In two subsequent steps the Illumina sequencing overhang adapters are introduced, which
38 then allow readout of the output. Representative gel images of the different stages can be found in
39 Supplementary Fig. 4.

40

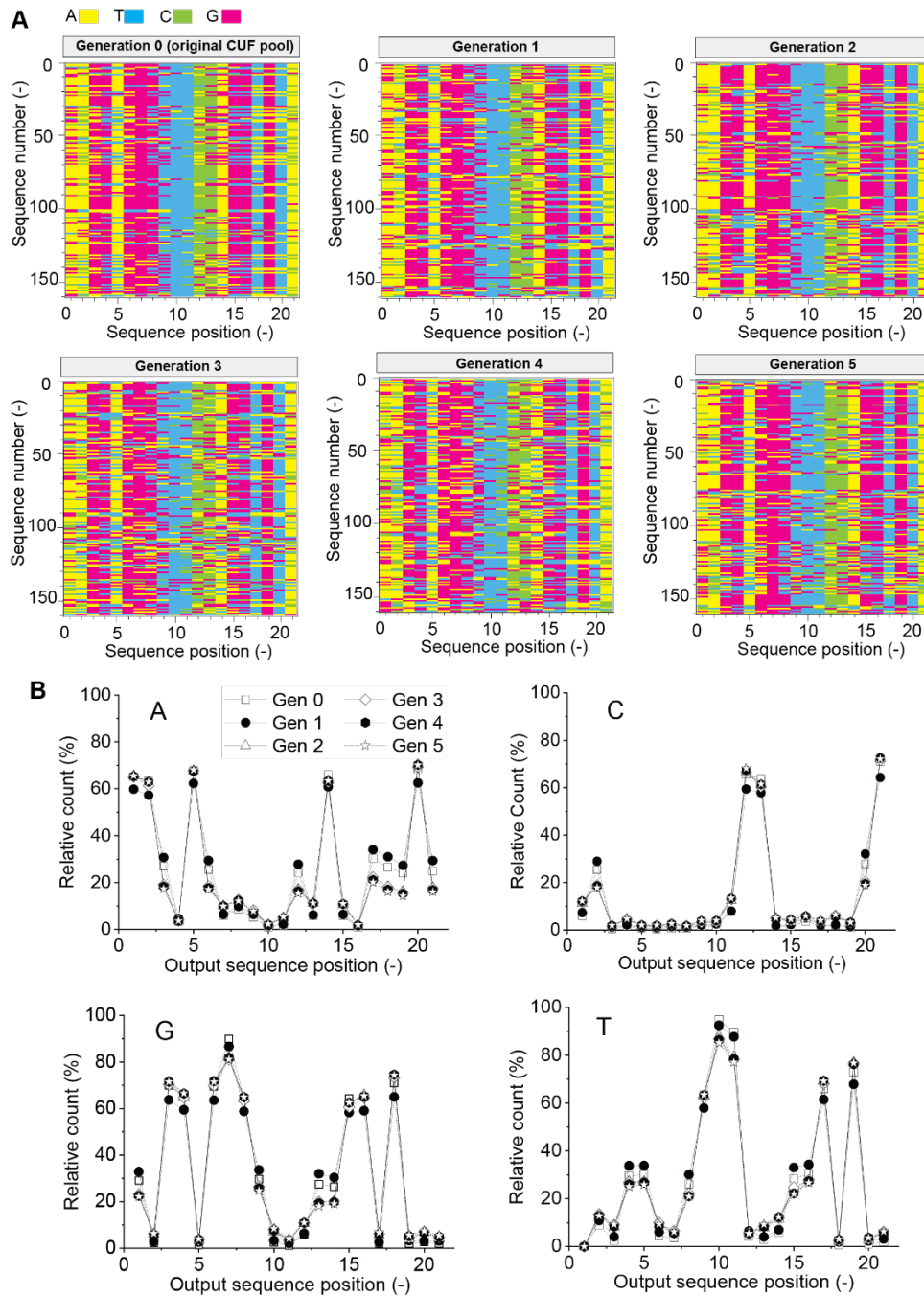


41

42 **Supplementary Fig. 2: Library design with cleavable handles.**

43 Shown are the sequence structures at various stages of the process of creating and operating the
 44 chemical random function. The working principle is identical as described for the library described in
 45 Supplementary Fig. 1, but there are additional steps to make the orDNA not only operable, but
 46 unclonable. The double-stranded (ds) PCR pool is treated with a type IIS restriction enzyme, PstI, that
 47 removes a large part of the handle and leaves a 5' overhang of a degenerate nucleotide. The sticky

48 end is blunted using Sequenase, which incorporates the complementary nucleotide with a dideoxy
49 modification.

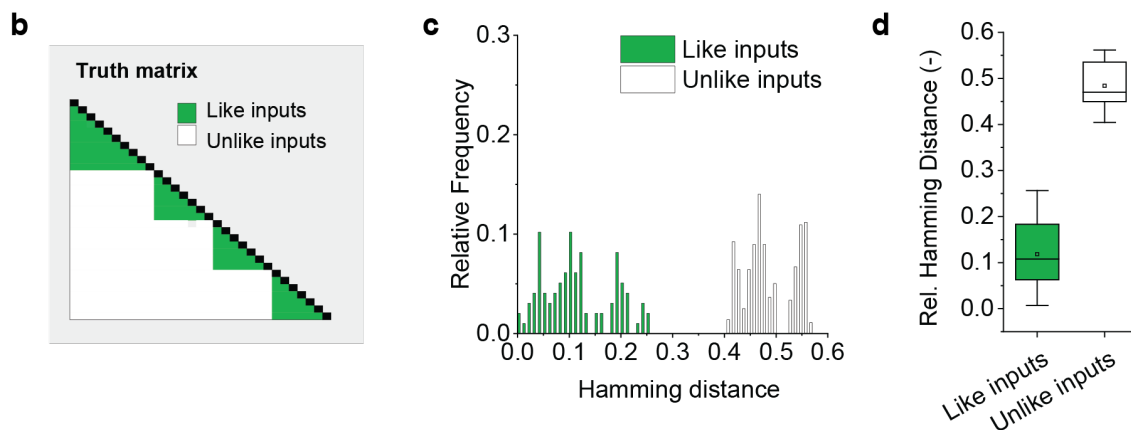
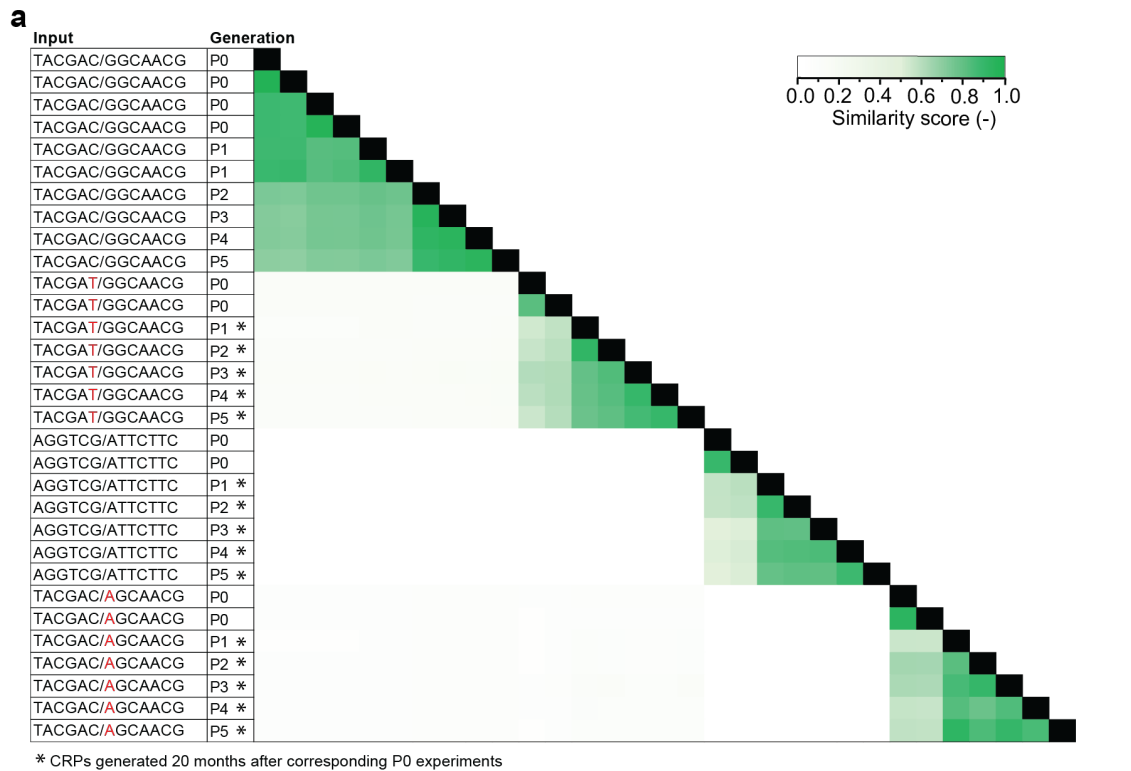


50

51 **Supplementary Fig. 3: Raw data analysis of proliferation CRPs**

52 Position-dependent raw data analysis of experiments 1 and 6-10, comprising the Illumina sequence
 53 reads resulting from the same input applied to five generations of the same CUF, i.e. the same orDNA
 54 pool, showing a high qualitative similarity between all generations (also refer to Fig. 3 (b), (c)). Each
 55 generation was created by copying the previous generation using PCR. A) Color-coded visualization of
 56 the occurrence of the four nucleobases A, C, G and T across the 21 positions of the output
 57 sequences. Shown are the first 160 reads out of the respective FASTQ file after filtering for the

58 presence of the constant adapter regions. The selection is thus arbitrary without any applied order. B)
59 Relative frequency of the four nucleobases A, C, G and T across the 21 positions of the output across
60 the entire sequence set resulting from Illumina sequencing. Source data are provided as a Source
61 Data file.



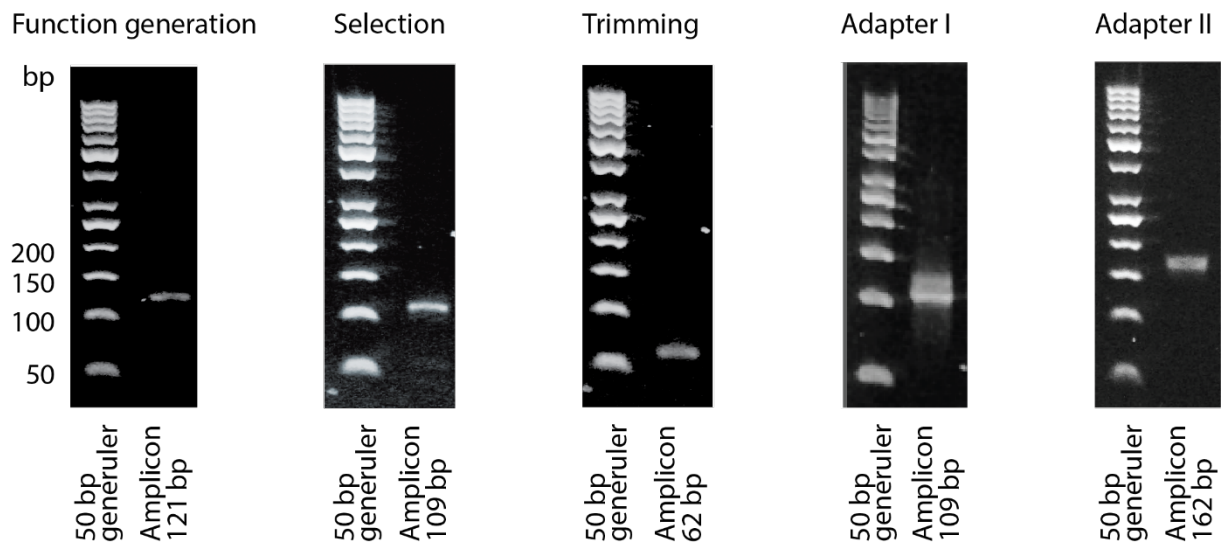
62

63 **Supplementary Fig. 4**

64 Extended data for generations P0-P5 (experiments 40-54, refer to Supplementary Tables 4-5), using
 65 four different inputs. Inputs 2-4 with generations P1-P5 were measured 20 months after the CRPs
 66 generated with P0. **a** shows the Jaccard similarity of CRPs comparing P0-P5 outputs. Jaccard
 67 similarity across generations as well as the 20-month time gap among like inputs is still well above
 68 Jaccard similarity of any CRP generated with unlike inputs, including the ones differing by a
 69 Levenshtein distance of 1. **b** Corresponding truth matrix to (a). **c** Histograms of Hamming distances
 70 after MinHashing showing the distributions of like and unlike inputs for experiments as shown in (a). n
 71 = 357 comparisons for unlike distribution, n = 98 for like distribution. **d** Boxplot showing average

72 relative Hamming distances, as per the distributions in (c). Indicated are the median (middle line),
73 mean (circled dot), 25th and 75th percentile (box) and 1.5 interquartile range (whiskers), with outliers
74 marked as black dots. n = 357 comparisons for unlike distribution, n = 98 for like distribution. Source
75 data are provided as a Source Data file.

76



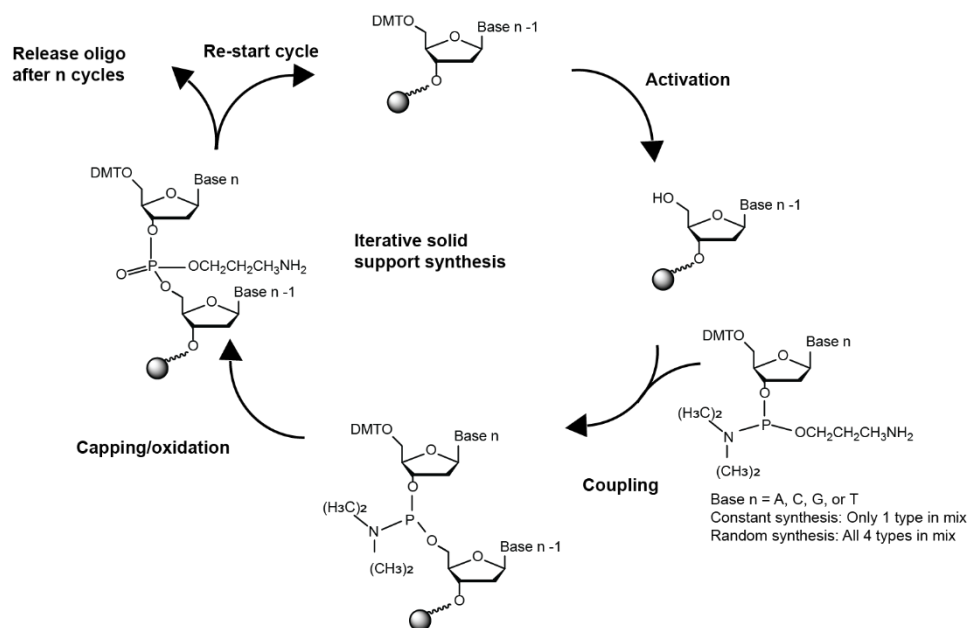
78

79 **Supplementary Fig. 5**

80 Agarose gel electrophoresis photographs of the different steps of function generation and operation,
 81 showing all stages of experiment 1. Images were converted from color to greyscale and edited in
 82 brightness and contrast using Adobe Photoshop.

83 **Supplementary Note 1: Random DNA synthesis**

84 Solid-support synthesis of DNA is commercially widely available. In standard chemical synthesis of
85 DNA, the strands are grown nucleotide by nucleotide in an iterative reaction, where each cycle
86 consists of several steps and adds a single base to each growing chain (Fig. 5). For most use cases,
87 this is done in a controlled manner, meaning only one type of building block (dNTP) at a time is added
88 to achieve billions of identical strands of constant length and known sequence. However, it is possible
89 to instead add mixtures of the four nucleobases A, C, G and T. In this case each growing chain will
90 incorporate only one of the four possible building blocks. It has been shown that this process is
91 random¹, meaning the entropy of a mixture can be exploited in a chemical reaction resulting in DNA
92 molecules of random sequence. Each sequence can thus be understood as a random (binary)
93 number, where every base contributes two bits. As the synthesis is stepwise, it can be individually
94 decided for each position along a synthesized sequence whether that specific position will contain a
95 determined or a randomly incorporated building block.



96

97 **Supplementary Fig. 6: Schematic representation of chemical DNA synthesis.**

98 The cyclic process adding a single nucleotide to each growing chain consists of an activation step,
99 followed by coupling, capping/oxidation. The cycle is followed until the desired chain length is
100 achieved, after which the synthesized oligonucleotide is released from the solid support. In the case of
101 sequence-determined synthesis, only a single type of nucleotide is added during the coupling step,

102 while for random synthesis, all four building blocks are mixed previous to addition, with each growing
103 chain incorporating only one of the four.

104

105 **Supplementary Note 2: Unpredictability and unknowability of sequence composition**

106 As in chemical unclonable functions each DNA strand contains segments that are randomly
107 synthesized, given the stochasticity of the process, it is a priori impossible to know the composition of
108 any given strand. For a chemical unclonable function based on operable random DNA this
109 unpredictability is key, as it translates to a random combination of input and output sequences.
110 Therefore, a given output sequence does not confer any information about the base composition of
111 the input and vice versa (see also Supplementary Note 3).

112

113 **Supplementary Note 3: Expected distribution of sequences within random pools and pool**
114 **entropy**

115 There are 4^n possible sequences for a randomly synthesized batch of DNA, where n corresponds to
116 the number of synthetic cycles using a dNTP mixture. The cycles using a single dNTP type, i.e. the
117 cycles for synthesis of the sequence-determined parts, are not included in the calculation, as they do
118 not contribute to the possibility space (aside from a negligible contribution by errors). As in this work,
119 libraries with a total of 40 random positions in the input and output segments were implemented, there
120 are $4^{40} = \text{ca } 10^{24}$ possible sequences, corresponding to approx. 1.66 mol of DNA. With an average of
121 330 g/mol/base · 121 bases (random and constant combined) this corresponds to a molecular weight of
122 ca. 40'000 g/mol for the single-stranded synthesis product of our library. This means to synthesize
123 every possible combination on average at least once, ca. 40'000 g/mol · 1.66 mol = 66 kg of DNA would
124 have to be produced. Around 4 mg of DNA were synthesized in total, equivalent to ca. $6 \cdot 10^{16}$
125 sequences. Importantly, the synthesis process only comprises individually growing chains without
126 copies (the sequences are only copied post-synthetically in a PCR reaction). This means any
127 duplicates only exist by chance and not by design. The probability for a specified pair of two individual
128 sequences being the same equals to $0.25^{40} = 8.27 \cdot 10^{-25}$.

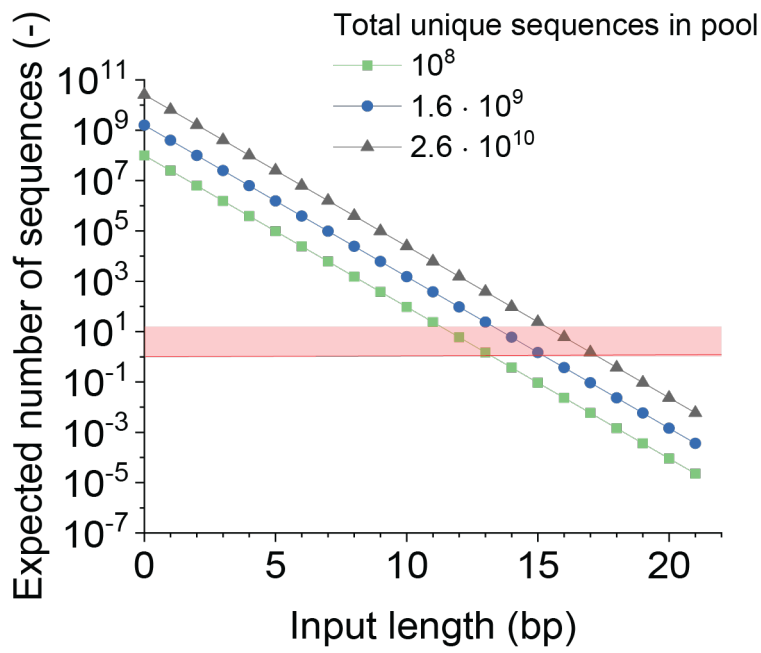
129 From the synthesized pool, 10^8 sequences, $1.6 \cdot 10^9$ and $2.6 \cdot 10^{10}$ sequences were arbitrarily extracted
130 to generate the three CUF sizes (with prefixes S, M and L). The extraction was performed by
131 dissolving the dried pool in ultrapure water to a known concentration. Using dilution series, it is then
132 possible to draw a volume that approximately corresponds to the number of desired sequences, which
133 are then subjected to PCR for function generation. With this ratio of function-determining vs. possible
134 combinations, the vast majority of the sequences is expected to be unique, with the overall average
135 copy number being very close to one.

136 As there are four possible building blocks per randomly synthesized sequence position, two bits of
137 entropy are encoded in a single base (translating to assigning binary numbers 00, 01, 10, 11 to the
138 four bases). As there are 40 random bases per sequence, this corresponds to 80 bits of random
139 information. As established above, ca 4 mg ($6 \cdot 10^{16}$ sequences or 100 nmol) of an orDNA library were
140 synthesized by a commercial supplier. This results in a total entropy of 80 bits/sequence $\times 6 \cdot 10^{16}$
141 sequences = $5 \cdot 10^{18}$ bits = 0.5 Exabyte (rounded to the first decimal position) in the pool.

142 **Supplementary Note 4: Design considerations for input primers**

143 Not all conceivable PCR primers are suitable input primers. Most importantly, the ratio between the
144 function size (= number of unique DNA sequences in a selection PCR, e.g. 10^8) and the chosen input
145 length in base pairs is relevant: If the input is too long, a given primer may not find any complementary
146 sequence in the pool to bind to, leading to an unspecific signal, or no amplification at all. If it is too
147 short, too many sequences may be amplified for successful signal extraction. This relationship is
148 shown in Fig. 6, which plots the expected number of output sequences against the input length for
149 three different pool sizes. For the pilot experiments, the input length was selected such that the
150 expected number of sequences that perfectly match to a given input primer sequence in the pool is
151 approx. 1.5. The actual number of matches varies between experiments and is expected to be
152 Poisson-distributed. For a pool of 10^8 sequences, the input length corresponding to this value is 13
153 random bases in total ($4^{13} = 6.7 \cdot 10^7$). The input was distributed over the forward and reverse primers
154 with 6 and 7 input bases, respectively. For the primers to reach high enough melting temperatures, 14
155 additional nucleotides were added, which were chosen to overlap with the respective outer handle
156 sequences (see Supplementary Fig. 1). This partial overlap with the constant regions of the orDNA is
157 additionally desirable, as this encourages binding at the intended positions, i.e. the random input
158 segments and not the random output segment. This results in amplicons of constant length and
159 ensures the amplified sequences contain the output segment, which is the chemical readout of the
160 function necessary to calculate a valid output.

161 However, inputs with a higher number of selective nucleotides have been shown to still generate
162 reproducible outputs in the test setting, even though no perfect counterpart is expected to be present
163 in the pool. This can be explained by the fact that there are always primer-template hybrids that are
164 thermodynamically more favored than others, and the corresponding sequences are thus favorably
165 amplified, even if the match is imperfect. As the readout of the function is generated by the entirety of
166 the amplified sequences and their frequency of occurrence, this distribution still works as a
167 reproducible signature. However, with an increased input length, a lower proportion of the nucleotides
168 is expected to contribute to the specificity of the readout and the noise increases, potentially leading to
169 the occurrence of collisions with shorter inputs. Therefore, depending on the application scenario, it is
170 advantageous to restrict the input length to conform to the respective pool size (see also
171 Supplementary Note 8.



172

173 **Supplementary Fig. 7**

174 Relationship of the expected number of matching output sequences in a pool related to the input
 175 length in basepairs (bp) at different orDNA pool diversities. The red area approximately marks the
 176 range this work operated in. Source data are provided as a Source Data file.

177

178 **Supplementary Note 5: Filtering of the sequencing reads**

179 Before applying k-mer extraction and computation of similarity, the reads were filtered to remove any
180 artefacts and potential contaminations. All correctly amplified orDNA sequences are expected to have
181 the same overall arrangement of their constant and random regions. Therefore, using the BBmap
182 (v38.99) command below, the reads were filtered for the presence of the constant primer region,
183 allowing a maximum Hamming distance of 3. Of the reads passing the first filter, those were excluded
184 that did not contain a 21-bp long insert between the expected constant segments. Only these inserts
185 of the passing reads, comprising the randomly synthesized 'output' portions, were included for further
186 analysis.

```
187 bbduk.sh in=<input.fq.gz> outm=stdout.fq ref=primer.fasta k=21 hdist=3 | bbduk.sh in=stdin.fq  
188 out=<output.fq.gz> ref=primer.fasta ktrim=r interleaved=f k=21 hdist=3 maxlength=21 mininsert=21  
189 minlength=21
```

190

191

192 **Supplementary Note 6: Choice of k-mer extraction and Jaccard similarity as data processing**
193 **methods**

194 In contrast to mathematical operations, data collected from real-world experiments are noisy. For the
195 assessment of chemical unclonable functions, this is an issue, as two individual operations of the
196 function will almost never lead to the exact same result, even if the equivalent inputs were used.
197 Physical unclonable functions and biometric data processing, e.g. human fingerprint analysis, face the
198 same issue. As a consequence, such datasets are analyzed in terms of their similarity instead of
199 perfect identity, with a defined similarity threshold above which the noisy datasets are considered
200 identical.

201 A k-mer-based method is suitable, as k-mers are commonly utilized in bioinformatics and, more
202 specifically, applied for genome comparisons². As the sequences of interest (the 'output' segments of
203 the sequences amplified during PCR selection with the input primers) are only 21 nt long, a circular k-
204 mer extraction was applied. This approach can be described by virtually circularizing the sequence of
205 interest and forming k-mers around the circle, as explained in Supplementary Fig. . For a sequence of
206 length n this leads to n k-mers independently of k , while linear extraction only yields $n-k+1$ k-mers.
207 Circular analysis therefore puts equal weight on all sequence positions, while a linear analysis would
208 over-emphasize the center nucleotides relative to the edges.

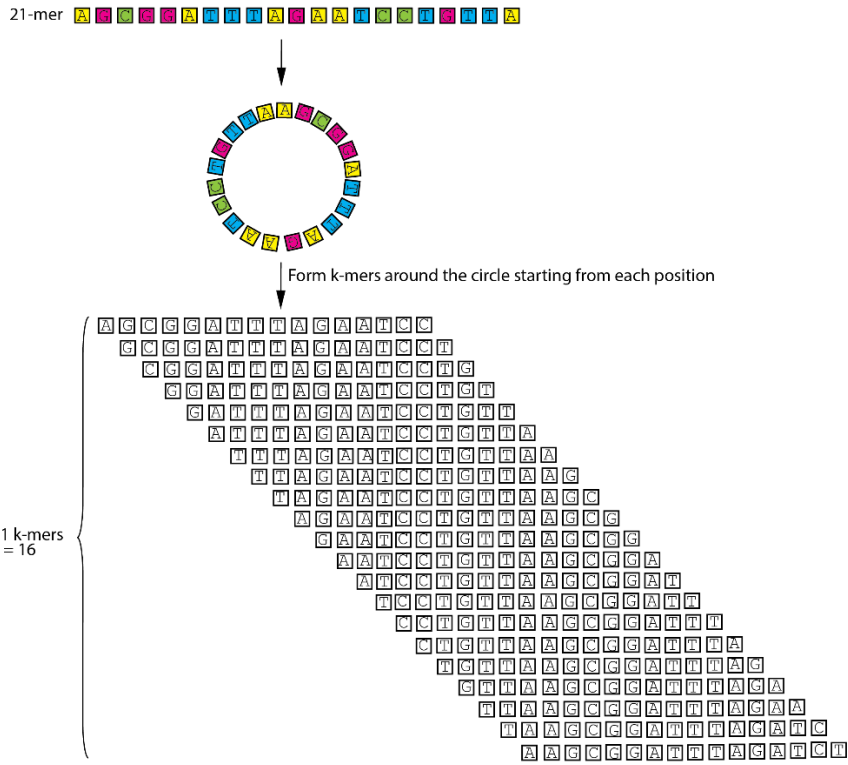
209 Further, to compare the k-mer sets extracted from the sequencing data, Jaccard similarity was used.
210 This index is widely used in machine learning, computational genomics and other fields³. It is
211 calculated by the ratio between the size of the intersection of two sets A and B divided by their union:

$$212 \quad J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

213 Furthermore, the Jaccard coefficient can be weighted:

$$214 \quad J(A, B) = \frac{\sum_k \min(A_k, B_k)}{\sum_k \max(A_k, B_k)}$$

215 While the unweighted coefficient only considers the presence or absence of a given k-mer, the
216 weighted coefficient considers the frequency of occurrence and thus the contribution of each k-mer to
217 the overall 'fingerprint'. Applying a weighted coefficient is therefore more comprehensive.



218

219 **Supplementary Fig. 8: k-mer analysis**

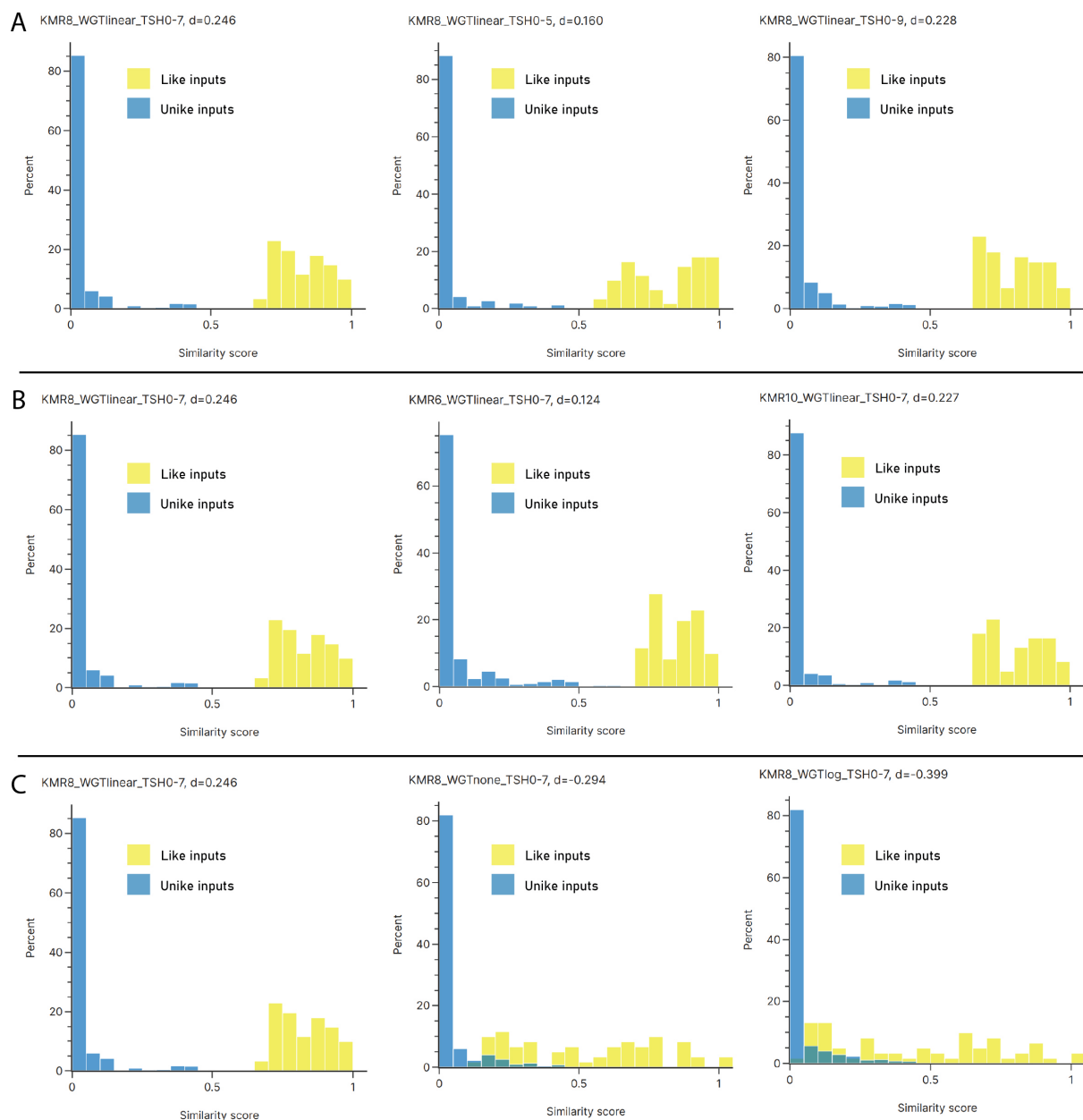
220 Schematic sketch of extracting all 21 possible 16-mers from a 21-mer in a circular k-mer extraction

221 procedure.

222 **Supplementary Note 7: Parameter choice and robustness**

223 The chosen methods of k-mer extraction and Jaccard similarity comprise several parameters that have
224 to be selected before data processing. Namely, k has to be selected, and optionally a weighting
225 regimen for the Jaccard coefficient as well as the choice of how many reads to include in the analysis.
226 The latter is designed to exclude the reads with the lowest counts from the analysis and is expressed
227 as a percentage threshold. I.e. a threshold of 0.7 only includes the reads that cumulatively make up
228 70% of the read count (starting with the highest frequency), and vice versa excludes the 30% lowest-
229 frequency reads. This is an additional measure to remove artefacts and noise while still utilizing the
230 relevant information.

231 These parameters influence the performance of the algorithm in correctly categorizing the compared
232 sets as belonging to like, or unlike inputs. In order to find well-performing parameters, they were
233 screened by combinatorically applying them to the collected original dataset of 39 individual
234 experiments (refer to Fig. 3(a), (c), and (d)). The screened values were $k = [4, 6, 8, 10, 12]$, weighting
235 = [none, linear, logarithmic] and read frequency threshold = [0.5, 0.7, 0.9]. The optimized metric was
236 the relative distance d of the Jaccard similarity score between the most similar sets resulting from
237 different inputs and the most dissimilar sets resulting from the same input. The parameters achieving
238 the highest distance between the two distributions (same and different inputs) were $k = 8$, linear
239 weighting and a read frequency threshold of 0.7 (70%). The separation between the like and unlike
240 input distributions is $d = 0.256$ when applying the selected parameters. The separation only slightly to
241 moderately decreases when changing the values of k and the threshold (Fig. 8a, b). However, the
242 separation is not maintained when switching the weighting regime (Fig. 8c). Reducing (by logarithm) or
243 removing the weighting does not significantly affect the similarity of unlike inputs but leads to a
244 broadening of the like input similarity distribution. This highlights that the read frequency information is
245 highly relevant for response generation, and that the differentiation of similar inputs is significantly
246 derived from changing frequencies rather than the mere presence of different k-mers.



247

248 **Supplementary Fig. 9: Parameter robustness**

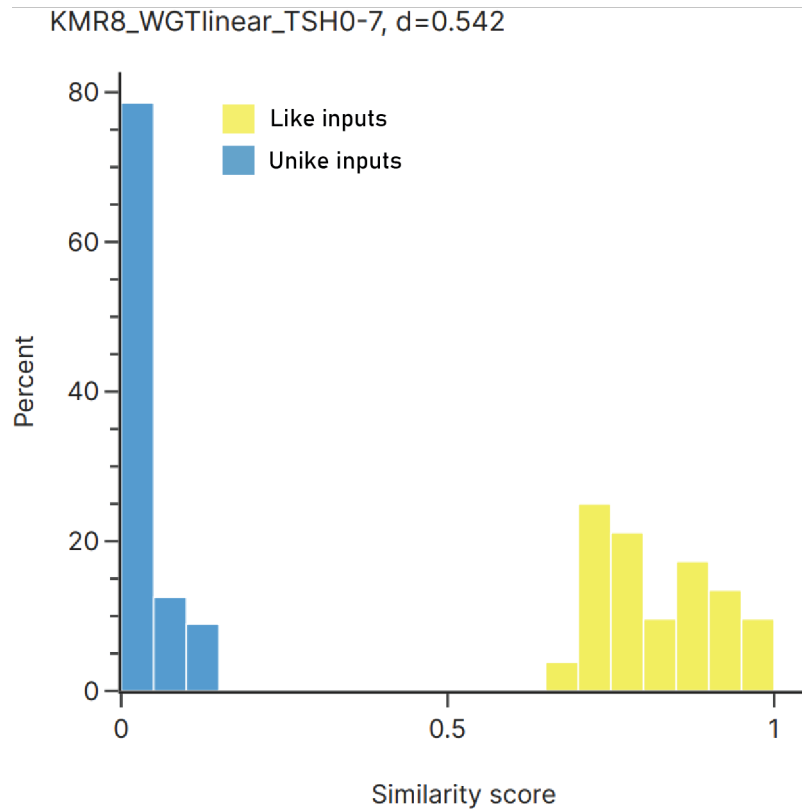
249 Histograms of similarity scores of all like (yellow distribution) and unlike (blue distribution) inputs
 250 across experiments 1-39 when using different parameters for k-mer extraction and similarity
 251 computation. The parameters and resulting distance d between the two distributions are indicated for
 252 each histogram, KMR referring to the value of k , WGT to the weighting applied (linear, logarithmic or
 253 none) and TSH to the relative read frequency threshold. A) Histograms showing the effect of varying
 254 the read frequency cut at $k=8$ with linear weighting. B) Histograms showing the effect of varying k
 255 while keeping the other parameters constant at TSH=0.7 and WGT=linear. C) Histograms showing the

256 effect of varying the weighting at constant $k=8$ and $TSH=0.7$. Source data are provided as a Source
257 Data file.

258 **Supplementary Note 8: Discussion of the parameters and input constraints**

259 It can be argued that a dataset comprising 39 experiments is not comprehensive and, therefore, that
260 the parameters emerging from the optimization are arbitrary. While this is a valid concern, the dataset
261 on which the parameters were optimized is heavily biased towards sets that are difficult to
262 differentiate. Many of the cross-compared inputs have a minimal Levenshtein distance of 1. Among
263 them were input primers that only differ in their length by a single nucleotide but were otherwise
264 identical, which are expected to draw a highly similar sequence set from the pool in the selection PCR.
265 The chosen parameters perform well in distinguishing these edge cases, showing an error rate of
266 zero. Moreover, they separate the like from the unlike input distribution by a large margin (the minimal
267 distance of Jaccard similarities being $d=0.256$, and the distance of the means $d=0.797$). Therefore,
268 these parameters can be expected to perform robustly over all potential inputs.

269 However, in a use case scenario, it may still be desirable to constrain the allowed inputs to a constant
270 primer length and to limit their GC-content for practical reasons. For example, this would allow running
271 all PCR reactions under the same conditions. When applying the length constraint in accordance with
272 the pool size (as discussed in Supplementary Note 4), the inputs are limited to 13 nucleotides in total,
273 distributed over the two primers. A common recommendation for optimal primer GC-content is
274 approximately 40-60%. Considering that in the current design the input primers contain a constant
275 portion with a fixed GC-content of 50%, the 13 nucleotides thus have to contain at least 3 and max. 10
276 G/C nucleotides. When applying these constraints to the existing dataset, the poly-T and the variable-
277 length inputs are excluded. The resulting distributions calculated with the previously selected
278 parameters ($k = 8$, linear weighting and a read frequency threshold of 0.7) are shown in Fig. 9. The
279 distance between the distributions approximately doubles from $d=0.256$ for the case without
280 constraints to $d=0.542$ when the length requirement and limited GC content are implemented. Notably,
281 this scenario still includes several inputs differing from each other by a Levenshtein distance of 1. This
282 further highlights the robustness of the chosen method of data processing.



283

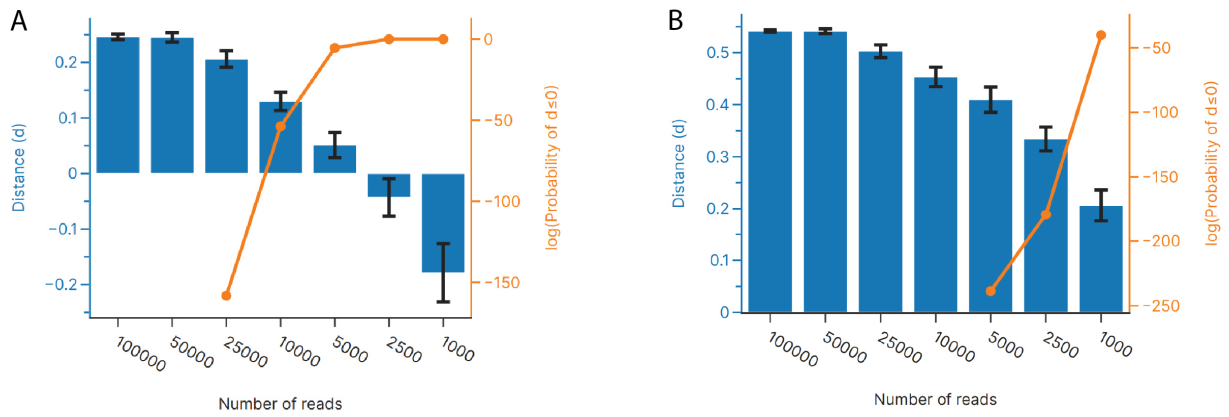
284 **Supplementary Fig. 10: Similarity scores with input constraints**

285 Histograms of similarity scores of all like (yellow distribution) and unlike (blue distribution) inputs
 286 across 24 experiments when applying input constraints regarding length and GC-content of the
 287 primers and using the parameters k=8, linear weighting and a read frequency threshold of 0.7. Source
 288 data are provided as a Source Data file.

289

290 **Supplementary Note 9: Down-sampling simulation**

291 There are experimental factors that lead to varying read numbers and quality per experiment. To
292 increase efficiency and extrapolate the limits of the system, down-sampling simulations were
293 conducted to find the average minimum number of reads per execution of the CUF that still allows for
294 correct output assignment with high fidelity. The simulation was performed with experiments 1-39
295 using a python pipeline and comprised stochastically drawing sets of different sizes out of the actual
296 number of reads using the python random number generator and performing feature extraction with
297 the drawn sets. The number of simulations per set size was 30. For samples with less reads than the
298 specified set size the entire read set was included in the analysis. The average distance of separation
299 between similarity scores of like and unlike outputs was then investigated with respect to the drawn set
300 size, the results of which are shown in Fig. 10. A distance of >0 indicates that the similarity
301 distributions of responses stemming from like and unlike inputs are entirely separable, meaning all
302 outputs can be correctly assigned. In contrast, a negative distance implies that certain edge cases
303 would be incorrectly assigned, i.e. two outputs stemming from the same input would be incorrectly
304 assigned as belonging to a different input, or vice versa. From the draws, the probability for one or
305 more response out of the generated dataset being incorrectly assigned was calculated with a one-
306 sided t-test, assuming a normal distribution. The simulations showed that below 10'000 reads, the
307 probability for a distance <0 becomes practically relevant in the case where variable-length inputs are
308 allowed. If length and GC-content constraints are applied (as discussed in Supplementary Note 9),
309 down-sampling indicates that even a read set comprising 1000 sequences can be assigned with a
310 success rate of close to 100%, even with the experiments including inputs with the minimal possible
311 distance from each other. This has a positive impact on affordability, as the sequencing cost
312 associated with 1000 reads is in the range of 10 cents.



313

314 **Supplementary Fig. 11: Distance (d) after downsampling to the respective number of reads.**

315 The distance refers to the minimal difference in similarity score between the sets resulting from like
 316 and unlike inputs after k-mer analysis using the entire dataset with parameters $k=8$, $\text{weight}=\text{linear}$,
 317 $\text{threshold}=0.7$. The number of draws per read was $n=30$. The error bars represent the 95% confidence
 318 intervals for the distance resulting from all draws. The curve shows the probability that at a given
 319 number of reads the resulting set distance using the indicated parameters will be ≤ 0 , whereby a
 320 negative distance means that like and unlike inputs can no longer be completely distinguished. For
 321 read numbers lacking a probability indicator the probability exceeded the float limit and thus shows as
 322 zero, i.e. is not represented in the log chart. A) Distance d after downsampling to the respective
 323 number of reads without applying input constraints. B) Distance d after downsampling to the
 324 respective number of reads in the case where the input is constrained with respect to GC-content and
 325 length. Source data are provided as a Source Data file.

326 **Supplementary Note 10: Use of MinHash signatures and fuzzy extraction for key generation**

327 A CUF outputs a set S of DNA sequences as a response to an input primer. If the function is evaluated
328 again with the same input primer, it outputs a set of DNA sequences that is similar to the previously
329 generated set S , and if the function is evaluated again with a different input primer, it outputs a set that
330 is different from S . In order to map a response sequence set S to an output of bits (or bytes), we use a
331 MinHash signature^{4, 5} followed by a fuzzy vault system⁶, which is a variant of a fuzzy extractor⁷.

332 MinHash-based tools are widely used, with numerous applications in genomics, such as taxonomic
333 diversity assessment⁸ and metagenome distance estimation⁹. A MinHash maps a set – in this case the
334 k -mers and associated abundances that comprise a response of a CUF – to a signature in form of a
335 vector. The vector dimensions are pre-defined and constant, irrespective of the set size, and the
336 similarity of two compared vectors approximately matches the Jaccard similarity of the respective k -
337 mer sets, from which the signatures were derived.

338 Specifically, our MinHash function uses a weighted MinHash implementation (datasketch Python
339 library) with the relative abundances of k -mers as weights. This, together with the same selection of k -
340 mer size and threshold, most accurately represents the analysis based on weighted Jaccard similarity
341 of Supplementary Note 7. Each result of the 255 MinHash permutations is converted to one byte by
342 modulo 256, yielding a signature of 255 bytes.

343 The vectors are then converted to binary keys by the fuzzy vault. This works as follows: Let w be the
344 data from a noisy process - in our setup w is generated from the MinHash signatures obtained from
345 the CUF. A fuzzy vault system generates a random key c , and uses it to generate the helper data h
346 from the data w . The helper data is publicly stored and provides redundancy for error correction but
347 does not give any advantage for reconstructing the original data w or the key c on its own. The original
348 data w is not stored, instead the generated random key c is used for authentication. If new data w' is
349 collected the fuzzy vault system attempts to reconstruct the original key c from the new data w' and
350 the original helper data h . If the new data w' is close (i.e., has a high similarity) to the original data w ,
351 then the redundancy included in the helper data h suffices to recover the same key c from w' .

352 More specifically, based on experiments 1-39, we implemented a fuzzy vault system with a so-called
353 code offset construction⁶. Given the original data w , we generate a 256-bit key c (i.e., 32 bytes) and
354 encode it with a with a Reed-Solomon code (reedsolo Python library, $255-32=223$ bytes of

355 redundancy), which yields a byte string c^* of equal length to the MinHash. We then store the helper
356 information $h = w - c^*$, where subtraction is in the finite field the letters of the codeword are in (i.e.,
357 bytes in our setup).

358 At the reconstruction phase, we are given the helper information h and data w' , and we compute

359
$$c' = \text{decode}(w' - h) = \text{decode}(w' - w + c^*)$$

360 where `decode` is the Reed-Solomon decoder. We have that $c' = c$ if the data w' is sufficiently close to
361 the original data w , since then, the redundancy afforded by the Reed-Solomon code allows
362 reconstruction of $c = \text{decode}(w' - w + c^*)$. Otherwise, reconstruction of c will fail, and no key is
363 generated.

364 See the code supplement for an implementation of the fuzzy vault system, which correctly reconstructs
365 all code words in the dataset stemming from experiments 1-39 (refer to Fig. 3g).

366 **Supplementary Note 11: Discussion of potential inverse operation and brute force attacks**

367 The response of a CUF to a challenge is the result of the random chemical composition of the orDNA
368 pool. This is a sufficient condition for the CRPs unpredictability. However, for cryptographic security,
369 irreversibility in analogy to a one-way function is needed. This means that a publicly stored output (and
370 the helper data) should not reveal information that can be used to infer the corresponding input. By the
371 CUF's materiality and random design, this is intrinsically the case. This is even further strengthened by
372 the fact that the chemical response – i.e. the sequence reads – is masked by k-mer analysis,
373 MinHashing and Fuzzy extraction, which are computationally infeasible to invert. So not only can the
374 input not be inferred from the output, but it is with current means even infeasible to unambiguously
375 infer the chemical response from the numerical output. This leads to a combination of digital and
376 materially manifested security layers.

377 However, the possibility of brute force attacks is unavoidable for any one-way function, PUF, or CUF.
378 For our CUF, the most straightforward way to prevent successful guessing of CRPs is to make the
379 pool so large that trial and error to find a given CRP is infeasible within reasonable time and cost.
380 Already at the currently used size of $2.6 \cdot 10^{10}$ sequences, a successful attack to find a CRP by trial
381 and error is unlikely, and CUFs even larger than the ones applied in this work would be conceivable.
382 To read all CRPs with the previously calculated 10'000 reads and at a pool size of $2.6 \cdot 10^{10}$
383 sequences, one would need ca 260 trillion reads (the equivalent of reading roughly half a million
384 human genomes at a depth of 30). Using state-of-the-art high-throughput sequencing such as the
385 NovaSeq X platform by Illumina, the estimated cost would approach half a billion USD for sequencing
386 alone, not accounting for PCR, time and labor, infrastructure, and other costs.

387 Additionally, this would require that an adversary malignantly gains physical access to the entire pool
388 in large enough quantities and knows the orDNA's general design. In our practical implementation of
389 the non-copiable CUF tokens, brute force attacks are therefore already prevented by the fact that only
390 a limited number of operations can be performed on a given token.

391 **Supplementary Note 12: Discussion of 2',3'-dideoxy end modification to introduce**
392 **unclonability**

393 As discussed in the main manuscript, a major constraint for re-creation of a chemical unclonable
394 function is the massive number of sequences, resulting in exorbitant cost of reading and recreating the
395 function. However, should an adversary gain physical access to a given CUF, using the outer handle
396 such an adversary could not only sequence, but also copy the pool using PCR.

397 Unclonability in the sense of uncopiability is a therefore desirable additional feature in the application
398 of a chemical random function. In order to switch an orDNA pool from the clonable (i.e. copiable) to the
399 unclonable (i.e. uncopiable) state, the constant handles at either end of the sequences are removed
400 via restriction digest. In the chosen implementation, only very short handles of 6 and 7 nucleotides,
401 respectively, remain on both ends. These are too short to function as universal PCR primer binding
402 regions due to their low melting temperature - the optimal T_m of primers being between 55 and 60 °C
403 with a recommended length of 18-30 nucleotides¹⁰.

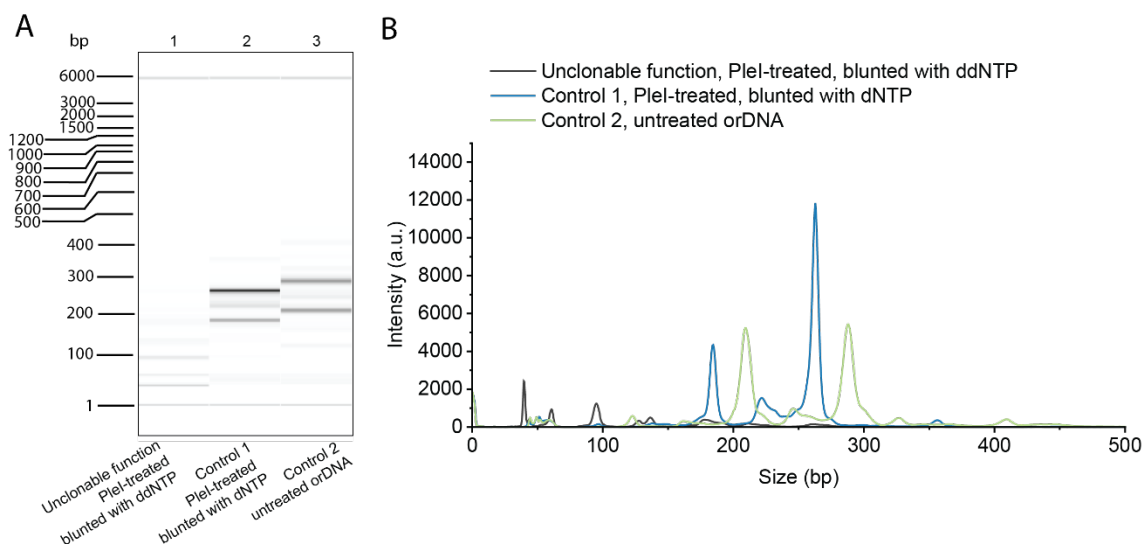
404 To prevent the re-addition of longer handles for PCR amplification, the sticky ends resulting from the
405 restriction digest are subsequently blunted using 2',3'-dideoxy nucleotides. Such dideoxy nucleotides
406 are used in Sanger sequencing¹¹ and their incorporation results in a 3' end without a free hydroxy
407 group. This means no ligation can be performed, as the 3'-hydroxy group is necessary for the formation of
408 a phosphodiester bond. The incorporation of dideoxy nucleotides is stable, i.e. cannot be efficiently
409 removed.¹²

410 To demonstrate that ligation is no longer possible, a CUF comprising approx. 10^8 unique sequences
411 treated with PstI and blunted with dideoxy nucleotides was sent to FASTERIS SA, an external service
412 provider, where it underwent a standard ligation procedure for genomic library preparation. The
413 selected ligation protocol is thus optimized for wide application and performs well in adding adapter
414 sequences to diverse libraries. Two controls were treated along with the test sample, one consisting of
415 the untreated orDNA library, the second one being digested and blunted, but using regular dNTPs
416 instead of dideoxy analogs. The total amount of DNA provided to the service provider was approx. 200
417 ng for all three samples, which is several orders of magnitude more than would be used for a PCR
418 reaction.

419 The ligation results and quality control data are shown in Supplementary Fig. . S11. While both control
420 libraries were successfully ligated with high yields, the dideoxy-treated CUF yielded no significant
421 product amounts of the expected length, as evidenced by the gel and the electropherogram. This
422 indicates that dideoxy blunting is indeed a suitable measure to make cloning of the function
423 impossible, or at least extremely difficult. Any intruder who attempts to copy such a CUF would not
424 only have to get physical access to an unrealistically high amount of the original function, but is likely
425 to introduce bias and change the structure of the orDNA in a way that would be evident and impair its
426 functionality.

427 The modification additionally introduces obstacles to simultaneous sequencing of the entire pool, as
428 this typically involves library preparation by adding adapters, either by PCR or by ligation. This also
429 includes “PCR-free” methods, such as nanopore sequencing. The truncated and modified 3'-ends
430 would need to be circumvented, which, even if possible, would increase the complexity and cost of
431 sequencing and copying further.

432 In conclusion, it can be shown that, aside from the necessity of gaining access to the pool with
433 malicious intent, there is no straightforward way to read and copy the pool. Any such attempt would
434 increase the effort, cost and complexity, and introduce the risk of altering the pool in a way that
435 hampers its performance as a random function.



436

437 **Supplementary Fig. 12: Ligation results**

438 A) gel image and B) electropherogram, as received from FASTERIS SA. The expected length of the
439 fully ligated fragments is 263 bp for samples 1 and 2, and 289 bp for sample 3. The fragments around
30

440 185 and 210 bp correspond to the one-ended ligation. The labels were added to the gel image by the
441 authors. The electropherograms were plotted from the raw data provided by FASTERIS SA. Source
442 data are provided as a Source Data file.

443

444 **Supplementary Note 13: Cost estimation**

445 The cost of a CUF is composed of the cost for the function itself and of its operation. For this study,
446 0.5 exabyte (4g of an orDNA library) were commercially ordered for a price of ca 90 USD. This
447 corresponds to less than 20 cents per 1000 Terabyte of entropy in the form of operable random DNA.
448 This amount of orDNA in theory suffices for billions of individual CUFs as implemented in this study,
449 meaning the cost of random synthesis is insignificant in any realistic scenario at scale. However,
450 generating a CUF from orDNA requires a PCR reaction to generate multiple copies, followed by
451 chemical modifications to switch the function from the clonable to the unclonable state. Within the
452 scales applied in this study, the costs for these steps approximately follow a linear relationship with
453 regards to the function size. For a CUF comprising 10^8 unique sequences, the singleplex qPCR using
454 an intercalating dye in 20 μ l reaction volume with standard primers costs ca 50 cents in reagents and
455 consumables at scale. The material generated by a single PCR reaction was on average sufficient for
456 ca. 200 executions of the CUF, but this is scalable.

457 The restriction digest with P1el costs ca 1 USD/assay in reagents and consumables, again for the
458 amount needed for 200 executions of the function. End blunting with ddNTPs costs ca. 3 USD;
459 although this calculation takes a prudent approach with the ddNTPs and enzyme added in large
460 excess. It is therefore likely that further process optimization would cut the cost significantly.

461 Adding up all cost items for function generation – from DNA synthesis to generating the double-
462 stranded orDNA to switching it to the unclonable state – results in a price of around 4.6 USD per CUF.

463 At least two executions are performed per input (one for registration, one or more for authentication
464 events). For a single function execution, 4 consecutive PCR reactions, including adding the
465 sequencing adapters, are performed, followed by Illumina sequencing. The cost for one operation from
466 input to output is comprised of reagents and consumables for PCR and next generation sequencing
467 (see Supplementary Table 1). For the sequencing cost, it was assumed that 10'000 reads are needed
468 for a single execution, which is a prudent estimate based on the down-sampling experiments
469 assuming no input constraints (as discussed in Supplementary Note 9). The current price of an iSeq
470 100 reagent kit was taken as a benchmark and broken down to the defined number of required reads.

471 In summary, the cost of generating the function lies below 1% and is thus insignificant. The combined
472 cost is instead dominated by the price for sequencing and PCR reagents, which currently make up

473 80% of the entire amount per function execution. It is important to note that the assumptions
474 underlying the sequencing cost estimate are deliberately conservative and based on the current
475 procedure; when reducing the required reads to 1000 and assuming the use of a large-scale
476 sequencing platform, the sequencing costs would drop by approx. two orders of magnitude. In that
477 case, the costs for PCR reagents and consumables would dominate the overall price at approx. 2
478 USD.

479 A caveat of the above calculations is that most of the items do not involve labor, as this is hard to
480 estimate in the setting of a research laboratory. Moreover, an assessment of the current labor
481 requirement would not be a realistic benchmark, as most of the operations are of low complexity and
482 would be automatized in a large-scale setting. The remaining hands-on work can easily be parallelized
483 in large multi-well plates, processing many function operations at the same time. Thus, at scale, labor
484 is not expected to have a large impact on the overall cost.

485

486 **Supplementary Table 1**487 Approximated item-by-item cost compilation for generation of a function comprising ca 10⁸ sequences.

Step	Cost item	Approximated cost at scale (USD)	Cost per execution*
Synthesis	Library	<0.01	<0.0001
orDNA generation	Mastermix (e.g. KAPA SYBR FAST, KAPA Biosystems, Wilmington, USA)	0.35	<0.01
	PCR primers	<0.01	<0.0001
	Consumables	0.15	0.0015
	Other reagents	0.10	0.001
	Restriction enzyme	0.80	0.008
End modification	Consumables	0.20	0.002
	Sequenase (e.g. ThermoFisher)	1.30	0.0013
	ddNTP mix	1.80	0.0018
	Other reagents	0.10	0.001
	Total		4.80

488 *calculated based on the assumption of 200 executions per function

489

490 **Supplementary Table 2**

491 Approximated item-by-item cost compilation for a single execution of a function comprising ca 10⁸
 492 sequences.

Step	Cost item	Approximated cost per function execution
Function operation	Mastermix (e.g. KAPA SYBR FAST, KAPA Biosystems, Wilmington, USA)	1.40
	PCR primers	0.05
	Consumables	0.60
	Illumina sequencing (iSeq 100 reagent kit)	1.50
	Other reagents	0.10
Total		3.65

493

494

495 **Supplementary Table 3: Sequence and primer list**

Name	Sequence 5'-3'
Library 1	ATGCGATGCAGTAAGCACTC NNNNNNNNNNNNNNNNNNNNNGCTCAGGATACCAAGCTGTCCNNNNNNNN NNGATATCTGCTCGGACCGCTA
Library 2	ATGCGAGTCAGATNGCACTC NNNNNNNNNNNNNNNNNNNNNGCTCAGGATACCAAGCTGTCCNNNNNNNN NNGACATNGGACGACTCAGCTA
Library 1 Handle primer fw	ATGCGATGCAGTAAGCACTC
Library 1 Handle primer rv	TAGCGGTCCGAGCAGATATC
Library 2 Handle primer fw	ATGCGAGTCAGATNGCACTC
Library 2 Handle primer rv	TAGCTGAGTCGTCCNATGTC
Input primer infw1	AGT AAG CAC TCG CTT ACG AC
Input primer invr1	GAG CAG ATA TCA TTG GCA ACG
Input primer infw4	TGCAGTAAGCACTCTACGAC
Input primer invr5	TCCGAGCAGATATCGGCAACG
Input primer infw4.1	TGCAGTAAGCACTCTACGAT
Input primer invr5.2	TCCGAGCAGATATCAGCAACG
Input primer infw4.3	TGCAGTAAGCACTCAGGTCG
Input primer invr5.3	TCCGAGCAGATATCATTCTTC
Input primer infw4.4	TGC AGT AAG CAC TCT TTT TT
Input primer invr5.4	TCC GAG CAG ATA TCT TTT TTT
Input primer infw5	GCAGTAAGCACTCTTACGAC
Input primer invr6	CCGAGCAGATATCTGGCAACG
Input primer infw6	CAGTAAGCACTCGTTACGAC

Input primer inv7	CGAGCAGATATCTTGGCAACG
Input primer inv8	GTCCGAGCAGATATCGCAACG
Input primer infw9	ATGCAGTAAGCACTCTACGA
Input primer inv9	ATGCAGTAAGCACTCTACGA
Trimming primer fw	ACACGACGCTCTTCCGATCT
Trimming primer rv	GGACAGCTTGGTATCCTGAGC
Illumina primer 1F	ACACTCTTTCCCTACACGACGCTCTTCCGATCT
Illumina primer 1R- AL	GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTGGACAGCTTGGTATCCT GAGC
Illumina primer 2FU	AATGATACGGCGACCACCGAGATCTACACTCTTCCCTACACGACGC
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATCGTGATGTGACTGGAGTTCAGACGTG
Index 1	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATACATCGGTGACTGGAGTTCAGACGTG
Index 2	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATGCCTAAGTGACTGGAGTTCAGACGTG
Index 3	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATTGGTCAGTGACTGGAGTTCAGACGTG
Index 4	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATCACTGTGTGACTGGAGTTCAGACGTG
Index 5	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATATTGGCGTGACTGGAGTTCAGACGTG
Index 6	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATGATCTGGTGACTGGAGTTCAGACGTG
Index 7	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATTCAAGTGACTGGAGTTCAGACGTG
Index 8	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATCTGATCGTGACTGGAGTTCAGACGTG
Index 9	T
Illumina primer 2RI	CAAGCAGAAGACGGCATAACGAGATAAGCTAGTGACTGGAGTTCAGACGTG
Index 10	T

Illumina primer 2RI Index 11	CAAGCAGAAGACGGCATAACGAGATGTAGCCGTGACTGGAGTTCAGACGTG T
Illumina primer 2RI Index 12	CAAGCAGAAGACGGCATAACGAGATTACAAGGTGACTGGAGTTCAGACGTG T
Illumina primer 2RI Index 13	CAAGCAGAAGACGGCATAACGAGATTTGACTGTGACTGGAGTTCAGACGTG T
Illumina primer 2RI Index 14	CAAGCAGAAGACGGCATAACGAGATGGAAGTGTGACTGGAGTTCAGACGTG T
Illumina primer 2RI Index 15	CAAGCAGAAGACGGCATAACGAGATTGACATGTGACTGGAGTTCAGACGTG T
Illumina primer 2RI Index 16	CAAGCAGAAGACGGCATAACGAGATGGACGGGTGACTGGAGTTCAGACGT GT
Illumina primer 2RI Index 17	CAAGCAGAAGACGGCATAACGAGATCTCTACGTGACTGGAGTTCAGACGTG T

496

497 **Supplementary Table 4: List of experiments**

Experiment number	Experiment designation	CUF description	Input primers	Input bases	Sequencing file name
1	1	CUF S1, 10 ⁸ sequences	infw4/inrv5	TACGAC / GGCAACG	Sp1-Infw4-Inrv5-1_S5_L001_R1_001.fastq
2	2	CUF S1, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	Sp1-Infw4-Inrv5-2_S6_L001_R1_001.fastq
3	3	CUF S1, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	Infw4-Inrv5-1_S1_L001_R1_001.fastq
4	4	CUF S1, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	Infw4-Inrv5-2_S2_L001_R1_001.fastq
5	P1.1	CUF S1 , Proliferation 1, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	Sp1-amp-Infw4-Inrv5-1_S11_L001_R1_001.fastq
6	P1.2	CUF S1 , Proliferation 1, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	Sp1-amp-Infw4-Inrv5-2_S12_L001_R1_001.fastq
7	P2	CUF S1 , Proliferation 2, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	0_S1_L001_R1_001.fastq
8	P3	CUF S1 , Proliferation 3, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	1_S2_L001_R1_001.fastq
9	P4	CUF S1 , Proliferation 4, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	2_S3_L001_R1_001.fastq

10	P5	CUF S1 , Proliferation 5, 10 ⁸ sequences	infw4/inrv5	TACGAC/ GGCAACG	3_S4_L001_R1_001.fastq
11	5	CUF S1, 10 ⁸ sequences	infw4.1/inrv5	TACGAT/ GGCAACG	Sp1-Infw4-1-Inrv5- 1_S9_L001_R1_001.fastq
12	6	CUF S1, 10 ⁸ sequences	infw4.1/inrv5	TACGAT/ GGCAACG	Sp1-Infw4-1-Inrv5- 2_S10_L001_R1_001.fastq
13	7	CUF S1, 10 ⁸ sequences	infw4/inrv5.2	TACGAC/ AGCAACG	Sp1-Infw4-Inrv5-2- 1_S11_L001_R1_001.fastq
14	8	CUF S1, 10 ⁸ sequences	infw4/inrv5.2	TACGAC/ AGCAACG	Sp1-Infw4-Inrv5-2- 2_S12_L001_R1_001.fastq
15	9	CUF S1, 10 ⁸ sequences	infw4.3/inrv5. 3	AGGTCG/ AATCATG	Sp1-Infw4-3-Inrv5-3- 2_S14_L001_R1_001.fastq
16	10	CUF S1, 10 ⁸ sequences	infw4.3/inrv5. 3	AGGTCG/ AATCATG	Sp1-Infw4-3-Inrv5-3- 2_S14_L001_R1_001.fastq
17	11	CUF S1, 10 ⁸ sequences	infw4.4/inrv5. 4	TTTTTT/ TTTTTTT	Lib4-Sp1-G2-infw4-4-inrv5- 4- 1_S13_L001_R1_001.fastq
18	12	CUF S1, 10 ⁸ sequences	infw4.4/inrv5. 4	TTTTTT/ TTTTTTT	Lib4-Sp1-G2-infw4-4-inrv5- 4- 2_S14_L001_R1_001.fastq
19	13	CUF S1, 10 ⁸ sequences	infw4/inrv8	TACGAC/ GCAACG	Sp1-Infw4-Inrv8- 1_S9_L001_R1_001.fastq
20	14	CUF S1, 10 ⁸ sequences	infw4/inrv8	TACGAC / GCAACG	Sp1-Infw4-Inrv8- 2_S10_L001_R1_001.fastq
21	15	CUF S1, 10 ⁸ sequences	infw9/inrv9	TACGA / GGCAAC	8_S9_L001_R1_001.fastq
22	16	CUF S1, 10 ⁸ sequences	infw9/inrv9	TACGA / GGCAAC	9_S10_L001_R1_001.fastq

23	17	CUF S1, 10 ⁸ sequences	infw1/inrv1	GCTTACGAC / ATTGGCAAC G	Infw1-Inrv1- 2_S3_L001_R1_001.fastq
24	18	CUF S1, 10 ⁸ sequences	infw1/inrv1	GCTTACGAC / ATTGGCAAC G	Sp1-Infw1-Inrv1- 1_S7_L001_R1_001.fastq
25	19	CUF S1, 10 ⁸ sequences	infw1/inrv1	GCTTACGAC / ATTGGCAAC G	Sp1-Infw1-Inrv1- 2_S8_L001_R1_001.fastq
26	20	CUF S2, 10 ⁸ sequences	infw4/inrv5	TACGAC / GGCAACG	Sp2-Infw4-Inrv5- 1_S1_L001_R1_001.fastq
27	21	CUF S2, 10 ⁸ sequences	infw4/inrv5	TACGAC / GGCAACG	Sp2-Infw4-Inrv5- 2_S2_L001_R1_001.fastq
28	22	CUF S2, 10 ⁸ sequences	infw1/inrv1	GCTTACGAC / ATTGGCAAC G	Sp2-Infw1-Inrv1- 1_S3_L001_R1_001.fastq
29	23	CUF S2, 10 ⁸ sequences	infw1/inrv1	GCTTACGAC / ATTGGCAAC G	Sp2-Infw1-Inrv1- 2_S4_L001_R1_001.fastq
30	24	CUF S3, 10 ⁸ sequences	infw4/inrv5	TACGAC / GGCAACG	Sp3-Infw4-Inrv5- 1_S3_L001_R1_001.fastq
31	25	CUF S3, 10 ⁸ sequences	infw4/inrv5	TACGAC / GGCAACG	Sp3-Infw4-Inrv5- 2_S4_L001_R1_001.fastq
32	26	CUF M1, 1.6 x 10 ⁹ sequences	infw5/inrv6	TTACGAC / TGCAACG	Sp-e9-Inrv5-Inrv6- 1_S5_L001_R1_001.fastq

33	27	CUF M1, 1.6 x 10 ⁹ sequences	infw5/inrv6	TTACGAC / TGGCAACG	Sp-e9-Inrv5-Inrv6- 2_S6_L001_R1_001.fastq
34	28	CUF M1, 1.6 x 10 ⁹ sequences	infw5/inrv5	TTACGAC / GGCAACG	Sp-e9-Inrv5-Inrv5- 1_S7_L001_R1_001.fastq
35	29	CUF M1, 1.6 x 10 ⁹ sequences	infw5/inrv5	TTACGAC / GGCAACG	Sp-e9-Inrv5-Inrv5- 2_S8_L001_R1_001.fastq
36	30	CUF L1, 2.6 x 10 ¹⁰ sequences	infw6/inrv7	GTTACGAC / TTGGCAACG	Sp-e10-Infw6-Inrv7- 1_S15_L001_R1_001.fastq
37	31	CUF L1, 2.6 x 10 ¹⁰ sequences	infw6/inrv7	GTTACGAC / TTGGCAACG	Sp-e10-Infw6-Inrv7- 2_S16_L001_R1_001.fastq
38	32	CUF L1, 2.6 x 10 ¹⁰ sequences	infw6/inrv6	GTTACGAC / TGGCAACG	Sp-e10-Infw6-Inrv6- 1_S13_L001_R1_001.fastq
39	33	CUF L1, 2.6 x 10 ¹⁰ sequences	infw6/inrv6	GTTACGAC / TGGCAACG	Sp-e10-Infw6-Inrv6- 2_S14_L001_R1_001.fastq
40	C2 P1*	CUF S1 , Proliferation 1, 10 ⁸ sequences	infw4.1/inrv5	TACGAT/ GGCAACG	P1-infw4- 1inrv5_S4_L001_R1_001.fastq
41	C2 P2*	CUF S1 , Proliferation 2, 10 ⁸ sequences	infw4.1/inrv5	TACGAT/ GGCAACG	P1-infw4- 1inrv5_S5_L001_R1_001.fastq
42	C2 P3*	CUF S1 , Proliferation 3, 10 ⁸ sequences	infw4.1/inrv5	TACGAT/ GGCAACG	P1-infw4- 1inrv5_S6_L001_R1_001.fastq
43	C2 P4*	CUF S1 , Proliferation 4, 10 ⁸ sequences	infw4.1/inrv5	TACGAT/ GGCAACG	P1-infw4- 1inrv5_S7_L001_R1_001.fastq

44	C2 P5*	CUF S1 , Proliferation 5, 10 ⁸ sequences	infw4.1/inrv5	TACGAT/ GGCAACG	P1-infw4- 1inrv5_S8_L001_R1_001.fastq
45	C4 P1*	CUF S1 , Proliferation 1, 10 ⁸ sequences	infw4.3/inrv5. 3	AGGTCG/ AATCATG	P1-infw4-3inrv5- 3_S14_L001_R1_001.fastq
46	C4 P2*	CUF S1 , Proliferation 2, 10 ⁸ sequences	infw4.3/inrv5. 3	AGGTCG/ AATCATG	P1-infw4-3inrv5- 3_S15_L001_R1_001.fastq
47	C4 P3*	CUF S1 , Proliferation 3, 10 ⁸ sequences	infw4.3/inrv5. 3	AGGTCG/ AATCATG	P1-infw4-3inrv5- 3_S16_L001_R1_001.fastq
48	C4 P4*	CUF S1 , Proliferation 4, 10 ⁸ sequences	infw4.3/inrv5. 3	AGGTCG/ AATCATG	P1-infw4-3inrv5- 3_S17_L001_R1_001.fastq
49	C4 P5*	CUF S1 , Proliferation 5, 10 ⁸ sequences	infw4.3/inrv5. 3	AGGTCG/ AATCATG	P1-infw4-3inrv5- 3_S18_L001_R1_001.fastq
50	C3 P1*	CUF S1 , Proliferation 1, 10 ⁸ sequences	infw4/inrv5.2	TACGAC/ AGCAACG	P1-infw4inrv5- 2_S9_L001_R1_001.fastq
51	C3 P2*	CUF S1 , Proliferation 2, 10 ⁸ sequences	infw4/inrv5.2	TACGAC/ AGCAACG	P1-infw4inrv5- 2_S10_L001_R1_001.fastq
52	C3 P3*	CUF S1 , Proliferation 3, 10 ⁸ sequences	infw4/inrv5.2	TACGAC/ AGCAACG	P1-infw4inrv5- 2_S11_L001_R1_001.fastq
53	C3 P4*	CUF S1 , Proliferation 4, 10 ⁸ sequences	infw4/inrv5.2	TACGAC/ AGCAACG	P1-infw4inrv5- 2_S12_L001_R1_001.fastq

54	C3 P5*	CUF S1 , Proliferation 5, 10 ⁸ sequences	infw4/inrv5.2	TACGAC/ AGCAACG	P1-infw4inrv5- 2_S13_L001_R1_001.fastq
----	-----------	---	---------------	--------------------	---

498 Experiment designation corresponds to the nomenclature used in the main manuscript.

499 **Supplementary Table 5: List of Selection PCR parameters and results**

Experiment number	Experiment designation	Input primers	Input bases	Selection PCR cycles	Ct	Annealing temp. (Ta)
1	1	infw4/inrv5	TACGAC / GGC AACG	29	17.54	62
2	2	infw4/inrv5	TACGAC / GGCAACG	29	17.54	62
3	3	infw4/inrv5	TACGAC / GGCAACG	29	17.43	62
4	4	infw4/inrv5	TACGAC / GGCAACG	29	17.51	62
5	P1.1	infw4/inrv5	TACGAC / GGCAACG	30	14.4	62
6	P1.2	infw4/inrv5	TACGAC / GGCAACG	30	14.08	62
7	P2	infw4/inrv5	TACGAC / GGCAACG	25	16.25	62
8	P3	infw4/inrv5	TACGAC / GGCAACG	25	16.49	62
9	P4	infw4/inrv5	TACGAC / GGCAACG	25	16.68	62
10	P5	infw4/inrv5	TACGAC / GGCAACG	25	15.72	62
11	5	infw4.1/inrv5	TACGAT / GGCAACG	50	21.12	62
12	6	infw4.1/inrv5	TACGAT / GGCAACG	50	21.89	62
13	7	infw4/inrv5.2	TACGAC / AGCAACG	50	20.04	62

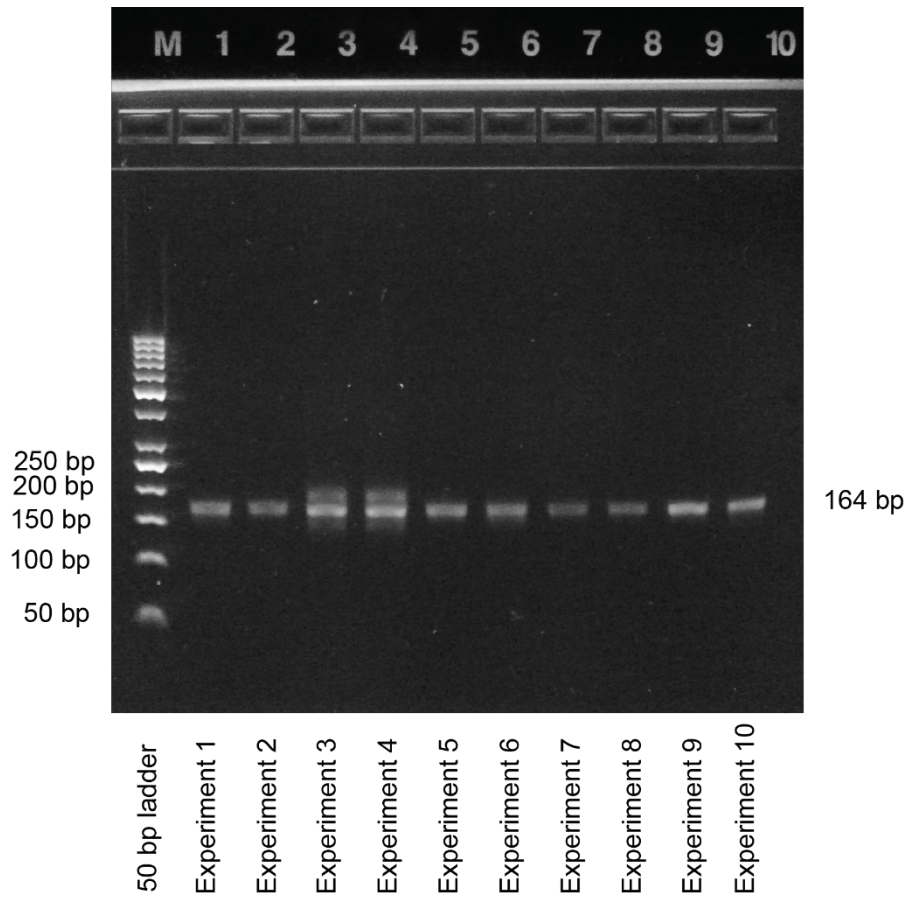
14	8	infw4/inrv5.2	TACGAC / AGCAACG	50	19.78	62
15	9	infw4.3/inrv5.3	AGGTCG / AATCATG	50	21.88	62
16	10	infw4.3/inrv5.3	AGGTCG / AATCATG	50	20.86	62
17	11	infw4.4/inrv5.4	TTTTTT / TTTTTTT	36	19.78	56
18	12	infw4.4/inrv5.4	TTTTTT / TTTTTTT	36	19.65	56
19	13	infw4/inrv8	TACGAC / GCAACG	25	15.5	62
20	14	infw4/inrv8	TACGAC / GCAACG	25	15.59	62
21	15	infw9/inrv9	TACGA / GGCAAC	24	11.9	62
22	16	infw9/inrv9	TACGA / GGCAAC	24	12.13	62
23	17	infw1/inrv1	GCTTACGAC / ATTGGCAACG	45	39.43	62
24	18	infw1/inrv1	GCTTACGAC / ATTGGCAACG	50	41.4	62
25	19	infw1/inrv1	GCTTACGAC / ATTGGCAACG	50	39.71	62
26	20	infw4/inrv5	TACGAC / GGCAACG	30	17.27	62
27	21	infw4/inrv5	TACGAC / GGCAACG	30	17.2	62
28	22	infw1/inrv1	GCTTACGAC / ATTGGCAACG	50	35.71	62

29	23	infw1/inrv1	GCTTACGAC / ATTGGCAACG	50	35.61	62
30	24	infw4/inrv5	TACGAC / GGCAACG	25	16.31	62
31	25	infw4/inrv5	TACGAC / GGCAACG	25	16.82	62
32	26	infw5/inrv6	TTACGAC / TGGCAACG	25	17.23	62
33	27	infw5/inrv6	TTACGAC / TGGCAACG	25	18.15	62
34	28	infw5/inrv5	TTACGAC / GGCAACG	25	17.35	62
35	29	infw5/inrv5	TTACGAC / GGCAACG	25	17.93	62
36	30	infw6/inrv7	GTTACGAC / TTGGCAACG	25	19.9	62
37	31	infw6/inrv7	GTTACGAC / TTGGCAACG	25	19.99	62
38	32	infw6/inrv6	GTTACGAC / TGGCAACG	25	18.69	62
39	33	infw6/inrv6	GTTACGAC / TGGCAACG	25	18.53	62
40	C2 P1*	infw4.1/inrv5	TACGAT/ GGCAACG	30	17.27	62
41	C2 P2*	infw4.1/inrv5	TACGAT/ GGCAACG	30	17.01	62
42	C2 P3*	infw4.1/inrv5	TACGAT/ GGCAACG	30	18.19	62

43	C2 P4*	infw4.1/inrv5	TACGAT/ GGCAACG	30	18.12	62
44	C2 P5*	infw4.1/inrv5	TACGAT/ GGCAACG	30	19.14	62
45	C4 P1*	infw4.3/inrv5.3	AGGTCG/ AATCATG	30	19.65	62
46	C4 P2*	infw4.3/inrv5.3	AGGTCG/ AATCATG	30	19.47	62
47	C4 P3*	infw4.3/inrv5.3	AGGTCG/ AATCATG	30	19.97	62
48	C4 P4*	infw4.3/inrv5.3	AGGTCG/ AATCATG	30	18.87	62
49	C4 P5*	infw4.3/inrv5.3	AGGTCG/ AATCATG	30	19.47	62
50	C3 P1*	infw4/inrv5.2	TACGAC/ AGCAACG	30	19.70	62
51	C3 P2*	infw4/inrv5.2	TACGAC/ AGCAACG	30	20.32	62
52	C3 P3*	infw4/inrv5.2	TACGAC/ AGCAACG	30	20.53	62
53	C3 P4*	infw4/inrv5.2	TACGAC/ AGCAACG	30	20.72	62
54	C3 P5*	infw4/inrv5.2	TACGAC/ AGCAACG	30	20.58	62

500

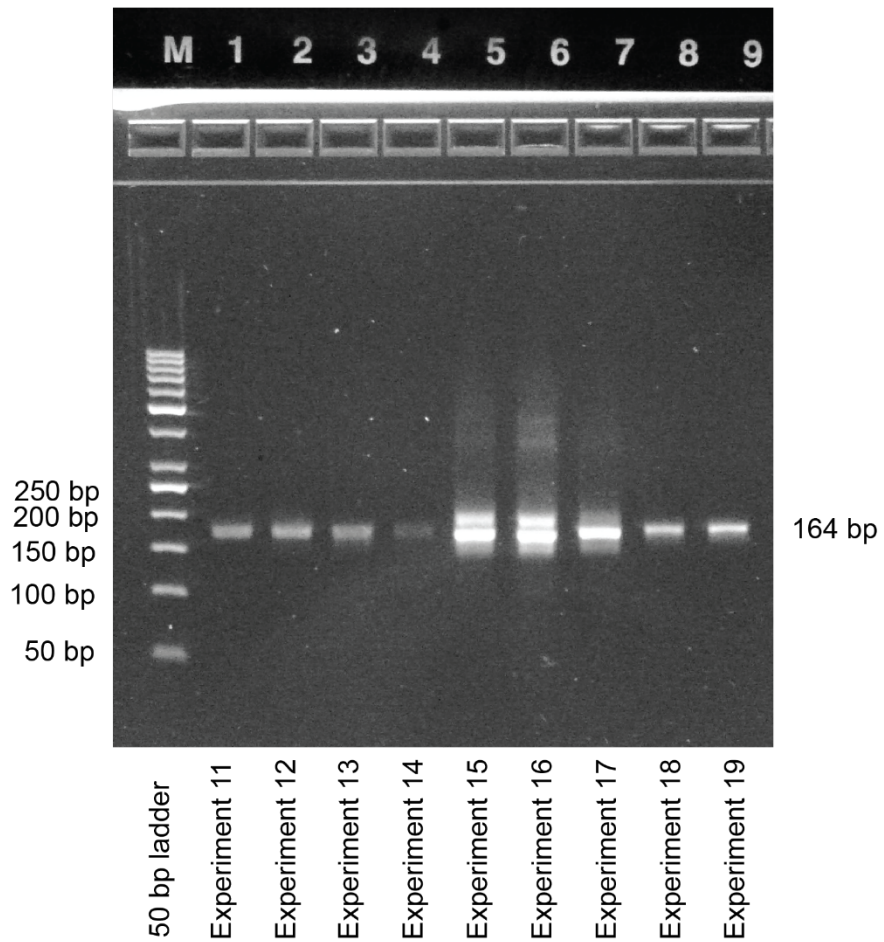
501



502

503 **Supplementary Fig. 13**

504 Agarose gel electrophoresis of experiments 1-10 (purified samples for sequencing). The image was
505 converted to grayscale and adjusted for brightness and contrast.

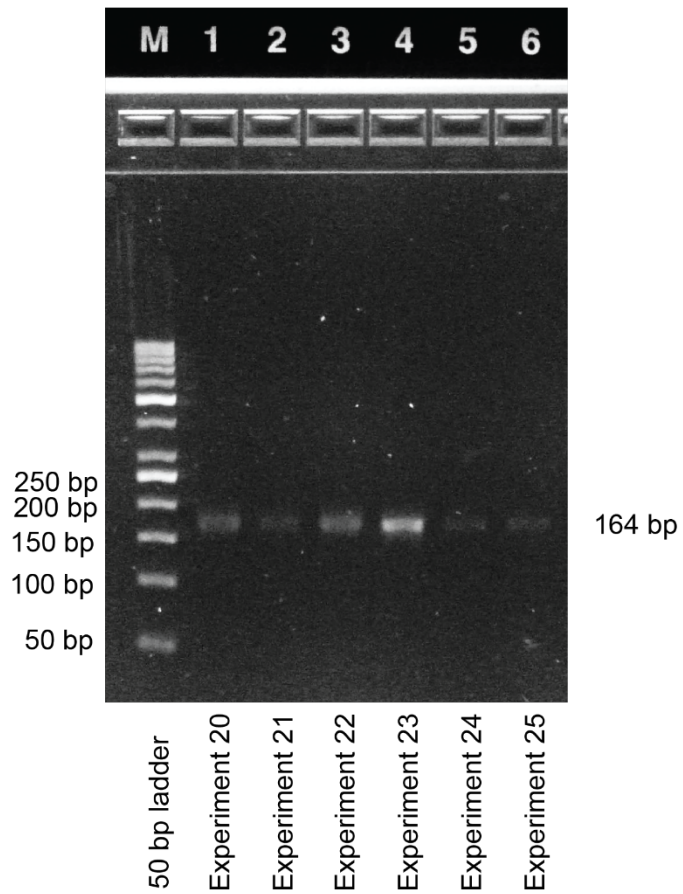


506

507 **Supplementary Fig. 14**

508 Agarose gel electrophoresis experiments 11-19 (purified samples for sequencing). The image was
509 converted to grayscale and adjusted for brightness and contrast.

510

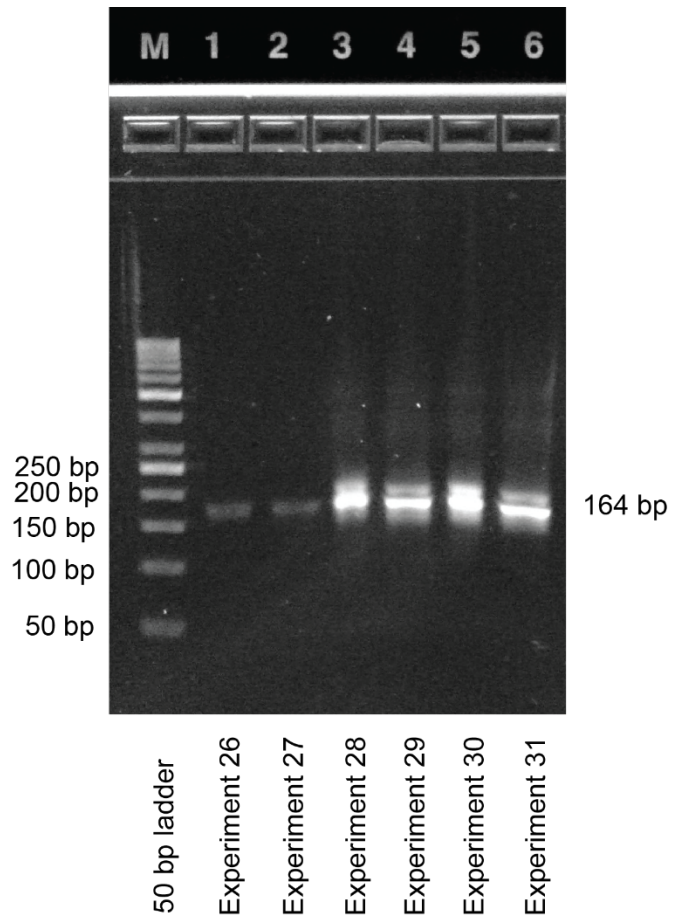


511

512 **Supplementary Fig. 15**

513 Agarose gel electrophoresis experiments 20-25 (purified samples for sequencing). The image was
514 converted to grayscale and adjusted for brightness and contrast.

515

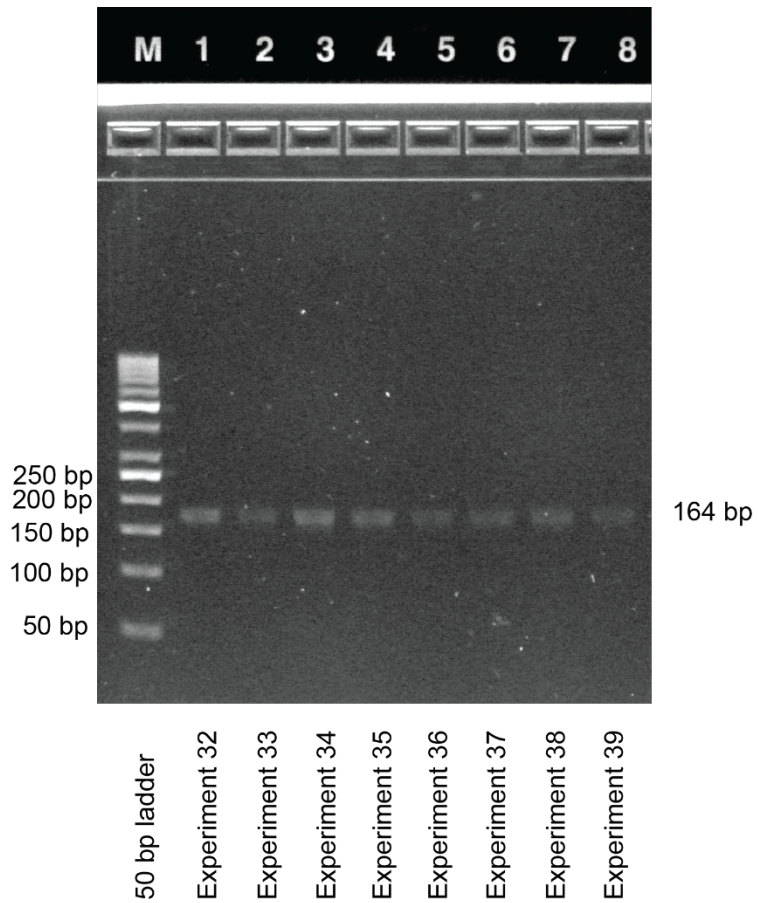


516

517 **Supplementary Fig. 16**

518 Agarose gel electrophoresis experiments 26-31 (purified samples for sequencing). The image was
519 converted to grayscale and adjusted for brightness and contrast.

520



521

522 **Supplementary Fig. 4**

523 Agarose gel electrophoresis experiments 32-39 (purified samples for sequencing). The image was
524 converted to grayscale and adjusted for brightness and contrast.

525

526 **References**

- 527 1. Meiser L. C., Koch J., Antkowiak P. L., Stark W. J., Heckel R. & Grass R. N. DNA synthesis for
528 true random number generation. *Nat. Commun.* **11**, 1-9 (2020).
- 529
530 2. Liu S. & Koslicki D. CMash: fast, multi-resolution estimation of k-mer-based Jaccard and
531 containment indices. *Bioinformatics* **38**, i28-i35 (2022).
- 532
533 3. Besta M., *et al.* Communication-efficient jaccard similarity for high-performance distributed
534 genome comparisons. In: *2020 IEEE International Parallel and Distributed Processing
535 Symposium (IPDPS)*. IEEE (2020).
- 536
537 4. Broder A. Z. On the resemblance and containment of documents. In: *Proceedings.
538 Compression and Complexity of SEQUENCES 1997 (Cat. No. 97TB100171)*. IEEE (1997).
- 539
540 5. Haveliwala T. H., Gionis A. & Indyk P. Scalable Techniques for Clustering the Web (extended
541 Abstract). In: *WebDB (informal proceedings)*. Citeseer (2000).
- 542
543 6. Juels A. & Sudan M. A fuzzy vault scheme. *Designs, Codes and Cryptography* **38**, 237-257
544 (2006).
- 545
546 7. Dodis Y., Ostrovsky R., Reyzin L. & Smith A. Fuzzy Extractors: How to Generate Strong Keys
547 from Biometrics and Other Noisy Data. *SIAM Journal on Computing* **38**, 97-139 (2008).
- 548
549 8. Katz K. S., Shutov O., Lapoint R., Kimelman M., Brister J. R. & O'Sullivan C. STAT: a fast,
550 scalable, MinHash-based k-mer tool to assess Sequence Read Archive next-generation
551 sequence submissions. *Genome Biology* **22**, 1-15 (2021).
- 552
553 9. Ondov B. D., *et al.* Mash: fast genome and metagenome distance estimation using MinHash.
554 *Genome biology* **17**, 1-14 (2016).
- 555
556 10. Apte A. & Daniel S. PCR primer design. *Cold Spring Harbor Protocols* **2009**, pdb. ip65 (2009).
- 557
558 11. Sanger F., Nicklen S. & Coulson A. R. DNA sequencing with chain-terminating inhibitors. *PNAS*
559 **74**, 5463-5467 (1977).
- 560
561 12. Lim S. E. & Copeland W. C. Differential incorporation and removal of antiviral
562 deoxynucleotides by human DNA polymerase γ . *Journal of Biological Chemistry* **276**, 23616-
563 23623 (2001).

564

565