

1 **Title**

2 Non-coding cancer driver candidates identified with a sample- and position-specific model of the
3 somatic mutation rate

4 **Authors and affiliations**

5 Malene Juul¹ (malene.juul.rasmussen@clin.au.dk)

6 Johanna Bertl¹ (johanna.bertl@clin.au.dk)

7 Qianyun Guo² (guo@cs.au.dk)

8 Morten Muhlig Nielsen¹ (morten.muhrig@clin.au.dk)

9 Michał Świtnicki¹ (michal.switnicki@clin.au.dk)

10 Henrik Hornshøj¹ (hhj@clin.au.dk)

11 Tobias Madsen¹ (tobias.madsen@clin.au.dk)

12 Asger Hobolth² (asger@birc.au.dk)

13 Jakob Skou Pedersen^{1,2} (jakob.skou@clin.au.dk)

14 ¹ Department of Molecular Medicine (MOMA), Aarhus University Hospital, Denmark

15 ² Bioinformatics Research Centre (BiRC), Aarhus University, Denmark

16 Corresponding authors: Jakob Skou Pedersen (jakob.skou@clin.au.dk) and Malene Juul

17 (malene.juul.rasmussen@clin.au.dk)

18 Abstract

19 Non-coding mutations may drive cancer development. Statistical detection of non-coding driver
20 regions is challenged by a varying mutation rate and uncertainty of functional impact. Here we
21 develop a statistically-founded non-coding driver-detection method, ncdDetect, which includes
22 sample-specific mutational signatures, long-range mutation rate variation, and position-specific
23 impact measures. Using ncdDetect, we screened non-coding regulatory regions of protein-
24 coding genes across a pan-cancer set of whole-genomes (n=505), which top-ranked known
25 drivers and identified new candidates. For individual candidates, presence of non-coding
26 mutations associate with altered expression or decreased patient survival across an
27 independent pan-cancer sample set (n=5,454). This includes an antigen-presenting gene
28 (*CD1A*), where 5'UTR mutations correlate significantly with decreased survival in melanoma.
29 Additionally, mutations in a base-excision-repair gene (*SMUG1*) correlate with a C-to-T
30 mutational-signature. Overall, we find that a rich model of mutational heterogeneity facilitates
31 non-coding driver identification and integrative analysis points to candidates of potential clinical
32 relevance.

33 Introduction

34 Cancer is caused by somatically acquired changes in the DNA sequence of genomes (Stratton,
35 Campbell, and Andrew Futreal 2009). Recently, large-scale sequencing of cancer-genomes
36 coordinated by the International Cancer Genome Consortium (ICGC), The Cancer Genome
37 Atlas (TCGA), and others has catalogued the molecular changes across hundreds of cancer
38 samples (International Cancer Genome Consortium et al. 2010; Cancer Genome Atlas
39 Research Network et al. 2013). The quest is now to analyze and understand the role of these
40 changes in cancer development. The aberrations in non-coding regions are of particular interest

41 as they have only become evident with the advent of whole cancer-genomes. Here we develop
42 the method ncdDetect for non-coding cancer driver detection. The method captures the
43 heterogeneities of the mutational process in cancer and aggregates signals of mutational
44 burden as well as functional impact in the significance evaluation of a candidate driver element.
45 We apply the method to 505 TCGA whole-genomes (Fredriksson et al. 2014).

46
47 Cancer arises by an evolutionary process where natural selection operates on genetic variation
48 stemming from randomly occurring somatic mutations. Thousands of somatic mutations
49 distinguish tumor tissue from healthy tissue, as a result of the mutational processes that cancer
50 cells go through during the lifetime of a cancer patient. The somatic mutations are identified
51 through Next Generation Sequencing (NGS) and commonly labelled according to their effect on
52 cancer development: *Driver mutations* are subject to positive selection during the evolutionary
53 process of cancer, as they offer the cell a growth advantage and contribute to the expansion of
54 tumors. By definition, driver genes contain one or more driver mutations. *Passenger mutations*,
55 on the other hand, have a neutral, or perhaps slightly negative, fitness contribution to the cell,
56 and accumulate as passive passengers during the evolutionary process of cancer (Stratton,
57 Campbell, and Andrew Futreal 2009; Pon and Marra 2015). Many more passenger than driver
58 mutations exist in cancer cells, and distinguishing the two is challenging (Marx 2014). Typically,
59 the strategy is to search for signs of recurrent positive selection across a set of cancer
60 genomes.

61
62 Signs of recurrent positive selection across cancer genomes can be detected by comparing the
63 somatic mutation frequency to an estimated background mutation rate (Pon and Marra 2015).
64 However, modelling the background mutation rate is complicated as it varies along the genome
65 with a large degree of heterogeneity (Lawrence et al. 2013; Polak et al. 2015; Alexandrov et al.
66 2013). Not only does the mutation rate in cancer exhibit high variation between different cancer

67 types; this is also the case between different samples of the same cancer type. Furthermore,
68 the mutational processes are affected by various genomic features, primarily replication timing,
69 expression, and the position-specific sequence context (Lawrence et al. 2013; Bertl et al., n.d.).
70 It is thus crucial to take these features into account when estimating the background mutation
71 rate in cancer. Another strategy for detecting signs of positive selection in cancer is to rank
72 mutations according to their impact on protein function (Marx 2014). In particular, point
73 mutations might introduce alterations in the amino acid sequence of a protein, and thereby
74 change its function (Reva, Antipin, and Sander 2011).

75
76 Systematic mutational screens of cancer exomes have expanded the set of known cancer driver
77 genes over the past decade (Forbes et al. 2014). Many tools exist for the identification of such
78 genes (Dees et al. 2012; Lawrence et al. 2013; Gonzalez-Perez and Lopez-Bigas 2012;
79 Tamborero, Gonzalez-Perez, and Lopez-Bigas 2013; Reimand and Bader 2013), and at present
80 616 cancer driver genes are catalogued as causally implicated in cancer (Cosmic 2016).
81 However, less than 2% of the genome codes for protein. While it is established that non-coding
82 elements play diverse roles in regulating the expression of protein-coding genes, only few
83 studies systematically explored the role of non-coding somatic mutations in cancer development
84 (Weinhold et al. 2014; Lochovsky et al. 2015; Melton et al. 2015; Fredriksson et al. 2014;
85 Mularoni et al. 2016; Lanzós et al. 2017). The first identified non-coding driver element was the
86 *TERT* promoter with highly recurrent mutations across several cancer types (Huang et al. 2013;
87 Horn et al. 2013). In general, the functional understanding of non-coding regions is poor
88 compared to protein-coding regions, challenging the interpretation of non-coding mutations
89 (Khurana et al. 2016).

90
91 We develop the method ncdDetect for detection of non-coding cancer driver elements. With this
92 method, we consider the frequency of mutations alongside their functional impact to reveal

93 signs of recurrent positive selection across cancer genomes. In particular, the observed
94 mutation frequency is compared to a sample- and position-specific background mutation rate,
95 which is estimated based on various genomic annotations. A scoring scheme (e.g. position-
96 specific evolutionary-conservation scores) is applied to further account for functional impact in
97 the significance evaluation of a candidate cancer driver element.

98

99 To strengthen our conclusions regarding the driver potential of candidate elements, we draw on
100 additional data sources. Non-coding mutations may perturb gene expression patterns, and we
101 thus correlate their presence with expression levels in an independent data set (Ding et al.
102 2015). Likewise, we correlate mutation status with survival information for these candidates.

103

104 What sets ncdDetect aside from other non-coding driver detection methods is the position-
105 specificity, and the derived ability to include genomic annotations of varying resolution down to
106 the level of individual positions. In one existing non-coding driver detection method, the position-
107 and sample-specific probabilities of mutation are derived, much like in ncdDetect, but are then
108 aggregated across a candidate element during significance evaluation (Melton et al. 2015). This
109 means that knowledge about the exact position and probability of a mutation is not fully utilized.
110 In another method, the genome is divided into bins according to the average value of replication
111 timing (Lochovsky et al. 2015), and in a recent method, the significance evaluation is performed
112 by locally conditioning on the number of observed mutations within a candidate element
113 (Mularoni et al. 2016). To our knowledge, no existing non-coding driver detection method derive
114 and apply position- and sample-specific probabilities of mutation in the significance evaluation of
115 a candidate driver element, and allows the use of position-specific scores and accurate
116 evaluation of their expectation across a candidate element. This unique feature of ncdDetect
117 means that candidate elements of arbitrary size and location can be analysed, and that the
118 potential large variation of mutational probabilities within a candidate element is handled.

119

120 With ncdDetect, we model the different levels of heterogeneity in the somatic mutation rate
121 known to be at play in cancer and evaluate the relative merit of different position-specific
122 scoring-schemes. The result is a driver detection method tailored for the non-coding part of the
123 genome, and with it we aim to contribute to the understanding of non-coding cancer driver
124 elements.

125 Results

126 ncdDetect evaluates if a given non-coding element is under recurrent positive selection across
127 cancer samples. The method takes as input (a) a candidate genomic region of interest, (b)
128 position- and sample-specific probabilities of mutation, and (c) position- and sample-specific
129 scores measuring mutational burden or impact.

130 Position- and sample-specific model of the background mutation rate

131 A key feature of ncdDetect is the application of position- and sample-specific probabilities of
132 mutation. These are obtained by a statistical null model, inferred from somatic mutation calls of
133 a collection of cancer samples (Material and methods: Statistical null model) (Bertl et al., n.d.).
134 The model predicts the mutation rate from a set of explanatory variables, i.e. genomic
135 annotations (Figure 1A). In the present paper, the null model is trained on 505 whole genomes
136 distributed across 14 different cancer types generated by TCGA (Fredriksson et al. 2014)
137 (Figure 1B).

138

139 As explanatory variables, the model includes genomic annotations known to correlate with the
140 mutation rate in cancer, as well as annotations we have found to improve the model fit. It is well
141 known that the mutation rate varies between samples (Lawrence et al. 2013; Alexandrov et al.

142 2013). Indeed, the mean and median number of mutations per sample is approximately $32 \times$
143 10^3 and 8×10^3 , respectively, with a large degree of variation between and within cancer types
144 (Figure 1B). For example, the average number of mutations per sample is 73 times higher for
145 colorectal cancer than for thyroid cancer, and within melanoma cancer, the least mutated
146 sample has 224 times fewer mutations compared to the highest mutated sample. The mutation
147 rate depends on the position-specific sequence context (Alexandrov et al. 2013) and correlates
148 with replication timing and gene expression level (Lawrence et al. 2013). The mutation rate also
149 varies between different types of genomic regions (Weinhold et al. 2014). Finally, we find that
150 the local mutation rate in a window around a given genomic position helps to capture
151 unaccounted mutation rate variation and increases the goodness of fit. The complete model
152 specification, including the definition of the location rate, is given in Material and methods:
153 Statistical null model. Consequently, genomic annotations considered as explanatory variables
154 in the null model for each sample are *replication timing*, *tissue-specific gene expression level*,
155 *trinucleotides* (the nucleotide under consideration and its left and right flanking bases, thus
156 taking into account the sample-specific mutational signature), *genomic segment* (3' and 5'
157 untranslated regions (UTRs), splice sites, promoter elements and protein-coding genes) and
158 *local mutation rate* (Figure 1C). Given these explanatory variables for a specific genomic
159 position, the model predicts the particular probability of observing a mutation of a given type for
160 a specific sample at this particular position.

161 Strand symmetric model

162 The reference sequence is divided into weak (A and T) and strong (G and C) base pairs (bps)
163 (Cornish-Bowden 1985). As strands generally cannot be distinguished in non-coding regions,
164 we handle them symmetrically, with weak bps denoted with an A and strong bps denoted with a

165 G. The 12 different types of point mutations are thus collapsed into the six classes, “A→C” (thus
166 including both A→C and T→G mutations), “A→G”, “A→T”, “G→T”, “G→C” and “G→A”.

167

168 Our model considers four possible outcomes for each position: Transitions ($TS_{\{A\rightarrow G, G\rightarrow A\}}$), two
169 types of transversions ($TV_{\{A\rightarrow T, G\rightarrow T\}}$ and $TV_{\{A\rightarrow C, G\rightarrow C\}}$) as well as the reference class of no
170 mutation.

171 Mutation-rate predictions and position-specific scores

172 For a given non-coding element of interest, the null model ensures the availability of position-
173 and sample-specific probabilities of each of the four possible outcomes (Figure 2A). Due to the
174 rarity of observing a mutation, the predicted mutational probabilities are of small magnitude
175 (Figure 2B-C, Figure 1-figure supplement 1). Additional to these probabilities, each position is
176 associated with a score that may be sample-specific, and may depend on the outcome class
177 (Figure 2D-E). The scoring scheme can be freely defined, and in the present paper we illustrate
178 three different choices of scores, namely the number of mutations, log-likelihoods and
179 conservation scores (Figure 2E). A wide variety of other scoring schemes can be considered. In
180 particular, the flexibility of ncdDetect allows for different scoring schemes to be chosen for
181 different types of candidate elements.

182 Scoring schemes

183 A good scoring scheme must be able to discriminate well between events that constitute true
184 cancer drivers and events that are neutral. The scoring scheme can, for example, evaluate the

185 mutational burden and be defined by means of the number of mutations in a candidate region.
186 This approach has been taken by existing non-coding cancer driver detection methods
187 (Lochovsky et al. 2015). Here, the score for a given position is defined to be one if a mutation of
188 any type occurs, and zero if no mutation occurs. Another approach is to evaluate the goodness
189 of fit of the observed mutations to the null model, and define the scores as log-likelihoods, i.e.
190 minus the natural logarithm of the predicted position- and sample-specific mutation probabilities.
191 This scoring scheme ensures that the more unlikely a mutational event, the higher the
192 associated score. A third approach is to also evaluate the functional impact of mutations when
193 defining the scores. However, for non-coding regions, we often lack the functional
194 understanding to interpret and predict the functional impact. We therefore illustrate this
195 approach using phyloP, a position-specific score of evolutionary conservation (Pollard et al.
196 2010), as a proxy score for functional impact (Material and methods: Scoring schemes).
197
198 The three proposed scoring schemes assign an *observed* score value of zero (number of
199 mutations and conservation scores), or a value close to zero (log-likelihoods), to positions with
200 no mutations, and a positive score to positions with mutations (figure 2E). The assigned scores
201 for mutated positions depend on the mutation type and the scoring scheme. Also, for each
202 position, all *possible* score values and the associated predicted probabilities are integrated in
203 the calculation of the background score distribution.

204 Driver detection

205 Though ncdDetect is designed for the analysis of non-coding elements, it can also be applied on
206 protein-coding genes. We initially evaluate the performance of different versions of ncdDetect
207 null models and different scoring schemes on protein-coding genes. We then use it to detect
208 driver candidates among promoter elements, splice sites, 5' UTRs and 3' UTRs (Table 1,

209 Material and methods: Candidate elements). While all of the analyses presented here are pan
210 cancer, individual cancer types can be analysed separately.

211 ncdDetect significance evaluation

212 With ncdDetect, significance evaluation of the observed mutations in a given genomic region of
213 interest is performed (Figure 3). ncdDetect uses a two-step algorithm in which sample-specific
214 calculations are followed by calculations across all samples in the dataset. The output is a p-
215 value indicating if the region of interest is under recurrent positive selection across the sample
216 set (Figure 3A).

217

218 The test statistic used in the significance evaluation performed by ncdDetect is the *observed*
219 *score*. This value is defined as the sum of sample- and position-specific observed scores across
220 the specific element that is being tested. For a given sample and a given position, the observed
221 score will depend on the chosen scoring scheme. For instance, the scoring scheme using the
222 number of mutations will always give a score of 1 to a mutated position, and a score of 0 to an
223 unmutated position. The scoring scheme using phyloP will give a score corresponding to the
224 position-specific phyloP value to a mutated position, and a score of 0 to an unmutated position.
225 Thus, the observed score for a specific element will depend on the chosen scoring scheme.

226

227 The observed score is significance evaluated in the *background score distribution*. Again, the
228 shape of this distribution will depend on the chosen scoring scheme: All possible sample- and
229 position-specific scores for the chosen scoring scheme are combined with the sample- and
230 position-specific mutational probabilities to form the background score-distribution.

231

232 The algorithm to obtain the background score-distribution works as follows: for a specific sample
233 *i*, the genomic region of interest is annotated with position-specific probabilities of mutation and

234 scores (Figure 3B-C). As noted above, the observed score of sample i is defined as the sum of
235 observed scores across all positions in the region (Figure 3D). In the current implementation,
236 only the two highest scoring mutations are considered for each sample. The position-specific
237 mutational probabilities and scores are aggregated using mathematical convolution; a method to
238 efficiently calculate the exact probability of observing a given sample-level score by summing up
239 the probabilities of all possible combinations of positional outcomes that could lead to it
240 (Grinstead and Snell 1997). The use of convolution is inspired by previously published protein-
241 coding driver detection methods (Lawrence et al. 2013; Dees et al. 2012). These calculations
242 lead to the sample-specific background score-distribution (Figure 3E). By repeating this
243 process, background score-distributions are found for each individual sample (Figure 3F).
244 These distributions are aggregated, again using convolution, to yield the overall background
245 score-distribution across samples. The individual sample-specific observed scores are added to
246 give the overall observed score, which is significance evaluated in the overall background
247 score-distribution. (Figure 3G) (more details are given in Appendix, section 3).

248 Protein-coding driver detection and model selection

249 To build a robust background null-model and evaluate the performance of ncdDetect, we apply
250 it to protein-coding genes (Figure 4). While we lack a well-established true-positive driver set for
251 the non-coding part of the genome, the COSMIC Cancer Gene Census (Cosmic 2016) provides
252 that for the protein-coding genes. As a performance measure, we therefore use the fraction of
253 COSMIC genes recalled among the ncdDetect candidate sets.

254

255 With ncdDetect, multiple hypothesis tests are performed. For example, protein-coding driver
256 detection requires significance evaluation of 19,256 genes. In order to evaluate all of these tests
257 simultaneously, QQ-plots are used to assess the distribution of the p-values and the number of
258 true hypotheses (Schweder and Spjøtvoll 1982). In these plots, the observed p-values are

259 plotted against the expected (uniform) p-values of the null distribution. P-values, which follow
260 the expected uniform distribution, will thus fall on the identity line, while smaller p-values will
261 deviate from this line. Per construction, 90% of the expected values lie in the interval $[1, 10^{-1}]$,
262 99% lie in the interval $[1, 10^{-2}]$, etc. (Figure 4A).

263 Model selection

264 The final model underlying ncdDetect is determined through a forward model-selection
265 procedure. In each step, position- and sample-specific probabilities are predicted for the protein-
266 coding genes, which are then evaluated with ncdDetect (Figure 4A, Figure 4-figure supplement
267 1). The *basic model* includes the genomic annotations sample id, replication timing and
268 trinucleotides as these are all known to correlate with mutation rate. The resulting p-values
269 appeared slightly inflated. To increase robustness of the predicted mutation-probabilities, we
270 defined *model 1a* by adding the variable local mutation rate to the basic model. This addition
271 resulted in less inflated p-values. However, the p-values were below the identity line in the QQ-
272 plot for more than 99% of the analysed genes, indicating that the predicted probabilities of
273 mutation were too large. As we found the somatic mutation rate is elevated in intergenic regions
274 compared to other functional elements, we defined *model 1b* by adding the variable genomic
275 segment to the basic model (Figure 1-figure supplement 1). This had the desired effect of
276 decreasing the final p-values, though leading to severe inflation. We defined *model 1c* that
277 extended the *basic model* with the local mutation rate, genomic segment, and tissue-specific
278 gene expression level. This lowered the p-values, although a small amount of inflation was still
279 observed.

280

281 Since we do not know all of the relevant genomic annotations that correlate with mutation rate
282 for all of our samples, it is unavoidable that we observe a difference between the actual and
283 predicted mutation rate (Figure A1). The effect of this difference will be accumulated along

284 elements, and even small biases in the predicted versus observed mutation rate may become
285 significant if elements are sufficiently long (Figure A2). The difference between the predicted
286 and observed mutation rate will cause overdispersion of the mutation rate. In the *final model*, we
287 thus correct for overdispersion by adjusting the sample- and position-specific probabilities of
288 mutation with an overdispersion-based rate adjustment (Materials and methods: An
289 overdispersion-based mutation rate adjustment, Appendix section 1). The resulting p-values
290 follow the expected uniform distribution, with less extreme p-values for the top-ranked genes
291 than for the previous models.

292 Recall of known protein-coding drivers

293 The p-values obtained with ncdDetect are corrected for multiple testing using a false discovery
294 rate of 10% (Benjamini and Hochberg 1995). The resulting ncdDetect candidate protein-coding
295 drivers are compared to the COSMIC Gene Census list for each of the three proposed scoring
296 schemes. We call 64 protein-coding genes significant using the conservation scores of which 15
297 (23%) are in COSMIC. In contrast, we call 109 protein-genes significant with log-likelihoods of
298 which 19 (17%) are in COSMIC, and 52 with the number of mutations of which 14 (27%) are in
299 COSMIC (Figure 4B, Supplementary files 1-3). The mean number of mutations per bp is on
300 average eight times higher for the COSMIC genes detected by ncdDetect compared to the
301 undetected COSMIC genes (Figure 4-figure supplement 2). The three proposed scoring
302 schemes have similar recall graphs, although the use of conservation scores appears the most
303 sensitive as it generally recalls the highest fraction of COSMIC genes (Figure 4C). For example,
304 in the top-15 protein-coding genes called by ncdDetect with conservation scores, nine are
305 COSMIC genes. This number is seven for the number of mutations, and seven for log-
306 likelihoods (Figure 4C, Figure 5-figure supplement 1A). The use of log-likelihoods results in the
307 highest number of elements called significant across most element types (Figure 5-figure
308 supplement 2).

309

310 As the mutational process is stochastic, it varies which drivers are involved in cancer
311 development, both within and between cancer types. COSMIC genes are identified from
312 analyses across many individual cancer types and a large fraction are likely not drivers in the
313 particular set of cancer samples analysed here. Furthermore, there might exist true protein-
314 coding cancer drivers not yet included in COSMIC. Out of the three proposed scoring schemes,
315 the conservation scores appear to have the highest sensitivity. It is more conservative than the
316 log-likelihoods, as it finds fewer significant protein-coding genes. Furthermore, it is compelling
317 that the use of this scoring scheme incorporates a measure of functional mutational impact in
318 the driver significance evaluation. In light of these considerations, we focus on the results
319 obtained with conservation scores in the following, and include the remaining two scoring
320 schemes for comparison.

321

322 To give an impression of how the calculated background score-distributions behave in practice,
323 we highlight a few examples (Figure 4D-F). The top-two protein-coding genes called by
324 ncdDetect are *TP53* (spanning 1,378 bps) and *PIK3CA* (spanning 3,207 bps), which are both
325 well-known cancer driver genes. An example of a protein-coding gene not called significant by
326 ncdDetect is *SLFN11* (spanning 2,706 bps). The smoothness of the overall score-distribution is
327 related to the length of the gene.

328

329 The performance of ncdDetect on protein-coding genes was benchmarked against a recent
330 non-coding cancer driver detection method, ExInAator (Lanzós et al. 2017) (Appendix, section 2).
331 In general, ExInAator predicts much fewer candidates than ncdDetect, and thus has a lower
332 false-positive rate. However, ncdDetect performs better at ranking genes compared to ExInAator
333 (Figure A3).

334

335 Non-coding driver detection

336 While the functional impact of non-coding mutations in cancer is not yet fully understood, it is
337 widely believed that they may play an important role in cancer development (Diederichs et al.
338 2016). Here, we apply ncdDetect to gene-associated non-coding elements of various types
339 (promoter elements, splice sites, 3' UTRs and 5' UTRs) to evaluate their cancer driver potential
340 (Figure 5, Figure 5-figure supplement 1, Supplementary Files 1-3, Material and methods:
341 Candidate elements).

342 Recall and function of previously described non-coding drivers

343 Promoter mutations might dysregulate gene expression in cancer. In particular, such mutations
344 might affect the expression levels of tumor suppressor genes or oncogenes (Diederichs et al.
345 2016). The average mutation rate in the analysed promoter elements is 7.0 mutations per mega
346 base (Mb) per sample (Figure 1-figure supplement 1). Of the investigated non-coding element
347 types, the promoter elements have the most significant calls in the ncdDetect analyses.
348 Approximately 1% of the evaluated promoter elements have a more significant p-value than
349 expected under the null (Figure 5A). We find a total of 160 significant ($q < 0.10$) promoter
350 elements. Within these, the observed mutation rate is 31.3 mutations per Mb per sample, which
351 is a fourfold increase of the mutation rate among all promoter elements. Of the promoter
352 elements, the *TERT* promoter is ranked most significant ($q = 2.4 \times 10^{-69}$, ncdDetect). The
353 promoter of the *TERT* gene is known to play an important role in telomerase expression, and
354 cancers with *TERT* promoter mutations have been shown to exhibit an elevated expression of
355 the *TERT* gene. This increased expression might ensure telomere maintenance, believed to
356 enable cancer cells to divide (Heidenreich et al. 2014). Two other identified promoter elements
357 are *WDR74* ($q = 4.1 \times 10^{-4}$, ncdDetect) and *PLEKHS1* ($q = 4.3 \times 10^{-5}$, ncdDetect). Mutations in
358 the promoter region of *WDR74* have been associated with increased gene expression and are

359 thought to have functional relevance for tumorigenesis (Khurana et al. 2013). Mutations in the
360 *PLEKHS1* promoter are also previously found to be significant in non-coding driver screens
361 (Weinhold et al. 2014; Melton et al. 2015). We note that out of the 863 whole genomes analysed
362 in (Weinhold et al. 2014), 356 are sequenced by TCGA. These samples appear to be a subset
363 of the 505 samples analysed in the present work, and the data sets are thus not completely
364 independent. In total, 29 of the 160 significant promoter elements called with ncdDetect are
365 previously found to be significant in non-coding cancer driver studies (Weinhold et al. 2014)
366 (Figure A4). As a benchmark of the performance of ncdDetect on regulatory non-coding regions,
367 we compared our results on promoter elements to those obtained with another non-coding
368 cancer driver detection method, LARVA (Lochovsky et al. 2015). The ncdDetect promoter
369 candidates that are not detected by LARVA include the previously described *WDR74*, *PLEKHS1*
370 and promoters of COSMIC genes (Appendix, section 2).

371

372 Another class of non-coding mutations are splice site mutations. They might disrupt the splicing
373 code and have been linked to cancer development (Srebrow and Kornblihtt 2006). The
374 destruction of a splice site will typically introduce stop codons or frameshifts and ruin the
375 function of the translated protein. The splice site mutation rate is 5.1 mutations per Mb per
376 sample. Three splice sites are found significant in this analysis (Figure 5B). As many as 90% of
377 the splice sites have zero observed mutations across the 505 samples. By construction, this
378 means that the resulting p-values are 1, and the p-value distribution is thus not uniform. More
379 samples would increase the detection power in these cases. Interestingly, in the top-ten ranking
380 splice sites we see a highly significant enrichment of splice sites associated to COSMIC genes
381 ($p = 6.1 \times 10^{-9}$, Fisher's exact test). Within the three significant splice site elements, the
382 mutation rate is 130.0 mutations per Mb per sample, corresponding to a 25-fold increase of the
383 mutation rate among all splice site elements. Splicing mutations in *TP53* are previously

384 described in cancer studies (Lee et al. 2010; Varley et al. 2001). Here, we observe that 12
385 samples from six different cancer types are mutated in the 52 bps that make up the splice sites
386 of *TP53*. These splicing mutations are highly significant ($q = 1.2 \times 10^{-21}$, ncdDetect), and might
387 lead to inactivation of the tumor suppressor *TP53* gene (Eicheler et al. 2002).

388
389 Finally, we investigate somatic mutations in the 3' and 5' UTRs (Figure 5-figure supplement 1B-
390 C), which regulate mRNA stability and translation. UTR mutations might disrupt binding sites for
391 miRNAs and RNA binding proteins and thereby affect post-transcriptional regulation. They might
392 also alter the structural conformations of the UTRs, which have previously been associated with
393 cancer (Diederichs et al. 2016). The average number of mutations is 6.4 per Mb per sample for
394 5' UTRs and 7.1 per Mb per sample for 3' UTRs. We find a total of 16 significant 3' UTRs and
395 86 significant 5' UTRs (Figure 5-figure supplement 1B-C). Within the significant 5' UTRs, the
396 mutation rate is 36.3 mutations per Mb per sample, a sixfold increase compared to all 5' UTRs.
397 For the significant 3' UTRs, the mutation rate is 20.6 mutations per Mb per sample, which is a
398 threefold increase of the average 3' UTR mutation rate. Two of the called 3' UTRs (*DRD5* and
399 *PCMTD1*) have previously been detected in cancer driver studies (Weinhold et al. 2014). This is
400 also the case for 12 of the 86 called 5' UTRs, one of them *SDHD* ($q = 2.3 \times 10^{-3}$, ncdDetect)
401 (Figure A4). In particular, a recent study identified the promoter region as well as the 5' UTR of
402 *SDHD* to be potential cancer drivers in melanoma. In particular, promoter mutations of *SDHD*
403 were shown to be associated with reduced gene expression and poor survival prognosis
404 (Weinhold et al. 2014). In the present data set, we observe six mutated melanoma samples in
405 the 5' UTR of *SDHD*, which covers 135 bps.

406 Case studies

407 The absence of a true-positive driver set for the non-coding part of the genome means that we
408 must find alternative ways to validate the driver potential of the candidates found by ncdDetect.
409 We thus seek to support the significance of the candidate elements and further characterize
410 them with evidence from two additional data sources.

411

412 A first approach is to analyse the effects of mutations on gene expression. To be able to look up
413 individual driver candidates, we gather expression values for each of the 505 considered whole
414 genome samples (Fredriksson et al. 2014) and perform a Wilcoxon rank sum test for top-
415 ranking candidates of each element type. To further support the findings, we obtain mutation
416 calls and expression values from the larger set of TCGA exomes (Cancer Genome Atlas
417 Research Network et al. 2013, Supplementary File 7) and likewise perform a rank sum test on
418 these data. (Supplementary File 5, Material and methods: Expression analysis). Reassuringly,
419 we recover known differences in *TERT* gene expression levels between samples mutated and
420 not mutated in the promoter region of the gene for the 505 whole genome samples ($q = 1.4 \times$
421 10^{-3} , Fisher's method) (Vinagre et al. 2013). Similarly for the TCGA exome samples, splice-site
422 mutations in *TP53*, which are known to drive cancer (Varley et al. 2001; Lee et al. 2010),
423 correlated with differences in gene expression levels ($q = 2.3 \times 10^{-2}$, Fisher's method).

424

425 Another approach we take to validate the ncdDetect candidates is to look at correlation between
426 mutation status and survival data for both the 505 whole genome samples and the TCGA
427 exomes. For this purpose, we download clinical data from the TCGA data portal ("TCGA Data
428 Portal" 2016). For a particular candidate driver, we test the significance of the difference in
429 survival between mutated and non-mutated samples using a one-sided Log-rank test on the

430 Kaplan-Meier estimated survival curves (Supplementary File 6, Material and methods: Survival
431 analysis). This recovers some known prognostic markers, such as *TP53* where splice site
432 mutations correlate with a significant decrease in survival ($q = 1.0 \times 10^{-1}$, Fisher's method)
433 (Yang et al. 2013). In the analysis of the 505 whole genome samples, we furthermore observe a
434 significant decrease in survival associated with *HLA-DRB1* promoter mutations ($q = 2.0 \times 10^{-2}$,
435 Fisher's method). Although this finding is potentially interesting, further investigation of this
436 candidate is beyond the scope of this paper, as genotyping in HLA regions is challenging due to
437 the highly polymorphic nature of these genes (Ehrenberg et al. 2014).

438

439 In the following we study a number of the top-ranking non-coding ncdDetect driver candidates in
440 detail. For each of them we further evaluate their driver potential by including results from
441 expression analysis and survival analysis.

442 *SMUG1 mutations and a uracil-DNA glycosylase deficiency mutational signature*

443 We observe 19 mutations in the 997 bp-long *SMUG1* promoter-region, which is approximately
444 seven times more than expected under the null model ($q = 1.1 \times 10^{-6}$, ncdDetect) (Figure 6A).

445 The mutations are distributed among 14 samples from three different cancer types (one breast
446 cancer sample, two colorectal cancer samples and eleven melanoma samples). As 15 out of 16
447 of the melanoma mutations are of type C→T in a CC context (or its reverse complement), they
448 are consistent with the mutational signature of ultraviolet (UV) light (Alexandrov and Stratton
449 2014). They may be a result of a mutational mechanism (Sabarinathan et al. 2016), however,
450 they may also affect *SMUG1* function.

451

452 *SMUG1* is involved in base excision repair (BER). Together with *UNG*, It acts in BER as an
453 uracil-DNA glycosylase, i.e., an enzyme that removes uracil from DNA (Visnes et al. 2009).
454 Uracil in DNA arises from spontaneous deamination of non-methylated cytosine, which causes
455 the occurrence of U:G mismatches. If unrepaired, they give rise to G→A transition mutations.

456 Mouse cell line experiments have shown additive effects of *SMUG1* and *UNG* inactivation on
457 G→A mutation rates (An et al. 2005). Furthermore, *UNG* and *SMUG1* expression has recently
458 been found to correlate negatively with genomic uracil levels in B cell lymphomas (Pettersen et
459 al. 2015). We therefore hypothesize that *SMUG1* mutations may affect the rate of G→A (and
460 C→T) mutations. To investigate this further, we define the *uracil-DNA glycosylase deficiency*
461 *signature* as the proportion of G→A (including C→T) mutations outside CpG sites (Figure 6B).

462 For melanoma samples, we further deduct C→T mutations in a CC context, as potentially
463 induced by UV light (Material and methods: Enrichment of G→A mutations in *SMUG1* mutated
464 samples).

465

466 There is a tendency for an increased value of the uracil-DNA glycosylase deficiency signature
467 statistic in *SMUG1* mutated melanoma samples ($p = 8.2 \times 10^{-2}$, one-sided Wilcoxon rank sum
468 test), and a significantly increased value of this statistic for the one *SMUG1* mutated uterus
469 cancer sample ($p = 2.1 \times 10^{-2}$, one-sided Wilcoxon rank sum test), compared to samples not

470 harbouring a *SMUG1* mutation (Figure 6C). This is also the case when restricting the analysis to
471 only include coding and splice site mutations (melanoma: $p = 5.3 \times 10^{-2}$, uterus cancer: $p = 2.1$
472 $\times 10^{-2}$, one-sided Wilcoxon rank sum tests). These findings indicate that *SMUG1* mutations
473 might perturb the uracil-DNA glycosylase function.

474

475 We further hypothesize that *SMUG1* and *UNG* expression may correlate with the uracil-DNA
476 glycosylase deficiency signature statistic. With the expression data available for the 505
477 analysed TCGA samples, we are unable to detect a significant correlation between gene
478 expression of *SMUG1* or *UNG* and the signature statistic (*SMUG1*: $p = 9.4 \times 10^{-1}$, *UNG*: $p = 1.0$
479 $\times 10^{-1}$, Fisher's method). To further investigate expression correlations, we look at the larger
480 data set of TCGA exomes. From these data, the uracil-DNA glycosylase deficiency signature
481 statistic is seen to be negatively correlated with *SMUG1* gene expression ($p = 3.8 \times 10^{-2}$,
482 Fisher's method), and with *UNG* gene expression ($p = 5.1 \times 10^{-4}$, Fisher's method). As *SMUG1*
483 and *UNG* are thought to have complementary roles in BER (Pettersen et al. 2007), we also
484 investigate the correlation between the signature statistic and the product of *SMUG1* and *UNG*
485 gene expression, which is negative and also significant ($p = 2.4 \times 10^{-3}$, Fisher's method) (Figure
486 6D, Figure 6-figure supplement 1).

487

488 Finally, we investigate whether survival correlates with *SMUG1* mutation status. With the
489 present data set, we are not able to detect such a pattern ($p = 8.2 \times 10^{-1}$, Fisher's method).

490

491 The observed correlations combined with the existing literature (Pettersen et al. 2015; An et al.
492 2005) suggest that mutations that functionally impact *SMUG1* and *UNG* may cause a mutational
493 phenotype as captured by the defined deficiency signature. However, further validation must
494 await availability of larger sets of cancer genomes.

495 *Promoter and UTR candidates where mutations associate with decreased survival*

496 The 5' UTR of *CD1A* spans 533 bp. In the region we observe 11 mutations from ten different
497 samples, distributed across eight different cancer types. This corresponds to five times the
498 amount of mutations expected under the null model ($q = 1.1 \times 10^{-2}$, ncdDetect). For TCGA
499 exome melanoma samples, we observe a highly significant decrease in survival associated with
500 mutations in the region ($p = 1.2 \times 10^{-9}$, Log-rank test) (Figure 7A), which is top-ranked among
501 the non-coding regions tested (Supplementary File 6). CD1 proteins present antigens to T cells
502 and are involved in eliciting adaptive immune responses. They are distantly related to HLA
503 (MHC) proteins and similarly bind T cell receptors, however, they display glycoproteins and
504 small molecules instead of peptides (Van Rhijn et al. 2015). Intriguingly, *CD1A* is generally lowly
505 expressed in healthy tissue with high expression particularly in skin (GTEx Consortium 2013),
506 where it is found in the antigen presenting Langerhans cells. *CD1A* has previously been
507 implicated with cancer development, with expression and positive correlation to survival
508 reported for some cancer types (Coventry and Heinzl 2004). Though we cannot functionally
509 interpret the observed TCGA melanoma mutations, the strong association with survival
510 suggests potential clinical relevance, not-the-least given the success of immunotherapy in
511 melanoma (Drake, Lipson, and Brahmer 2014).

512

513 A total of 22 mutations are observed in the 1,976 bp-long promoter region of *PRSS3*. This is
514 four times more mutations than expected under the null model, and they occur in 13 samples

515 from seven different cancer types ($q = 1.1 \times 10^{-2}$, ncdDetect). Previous studies have established
516 the role of *PRSS3* in the progression of pancreatic and ovarian cancer (Jiang et al. 2010; Ma et
517 al. 2015). Not only the promoter region of this gene comes out significant in the driver detection
518 screen; this is also the case for its 3' UTR ($q = 1.4 \times 10^{-5}$, ncdDetect) as well as its protein-
519 coding gene ($q = 6.8 \times 10^{-20}$, ncdDetect). Based on the TCGA exome set, we observe a
520 significantly decreased survival for head and neck cancer (HNSC) samples mutated in the
521 promoter region of *PRSS3* ($p = 1.8 \times 10^{-3}$, Log-rank test) (Figure 7B) as well as in the *PRSS3*
522 coding gene ($p = 1.2 \times 10^{-2}$, Log-rank test). We also observe a tendency for decreased survival
523 among melanoma samples with 3' UTR mutations ($p = 8.3 \times 10^{-2}$, Log-rank test).

524

525 The 3' UTR of *SEC14L1* spans 3,052 bps and contains 31 mutations from 27 samples
526 distributed across ten different cancer types. This is approximately four times the amount of
527 expected mutations ($q = 8.9 \times 10^{-5}$, ncdDetect), and the majority (58%) of these are found in
528 breast- and colorectal cancer. Although little is known about *SEC14L1* in cancer, one study
529 hypothesized that altered expression of the gene could contribute to breast tumorigenesis
530 (Kalikin et al. 2001). Another study found *SEC14L1* to be overexpressed in prostate cancer
531 (Burdelski et al. 2015). For TCGA exome HNSC samples, we find a significant decrease in
532 survival associated with mutations in the 3' UTR region of the gene ($p = 5.8 \times 10^{-5}$, Log-rank
533 test).

534 *STK11* splice sites mutations and their expression correlation

535 The combined splice site region of *STK11* covers 36 bps and is mutated in four lung
536 adenocarcinoma (LUAD) cancer samples, which is approximately 56 times more mutations than
537 expected under the null model (Figure 7C). ncdDetect ranks the splice site region of *STK11*
538 second among all splice sites ($q = 2.2 \times 10^{-3}$, ncdDetect). *STK11* is a known COSMIC tumour
539 suppressor gene, which has been shown to be involved in lung and cervical cancers (Gill et al.
540 2011), and very recently, splice site mutations of the gene were described in relation to cancer
541 (Mularoni et al. 2016; Wei et al. 2016). From the 505 whole genomes analysed here, we are
542 unable to associate the splice site mutations of *STK11* with a changed level of gene expression
543 ($p = 4.4 \times 10^{-1}$, Fisher's method). Looking at the larger set of TCGA exomes, however, we
544 detect a significantly lower expression level for LUAD samples mutated in the splice site region
545 of *STK11*, compared to LUAD samples without such mutations ($p = 1.3 \times 10^{-3}$, two-sided
546 Wilcoxon rank sum test) (Figure 7D). We further observe a marginally significant decrease in
547 survival associated with *STK11* splice site mutations for LUAD TCGA exome samples ($p = 6.5$
548 $\times 10^{-2}$, Log-rank test) (Figure 7E).

549 Discussion

550 Non-coding somatic mutations play part in tumour initiation and progression. With the advent of
551 whole genome sequencing, the systematic screening of such mutations is possible. We have
552 developed the method ncdDetect with the goal of detecting non-coding cancer driver elements
553 and thereby gain an understanding of the underlying mechanisms of tumorigenesis. With
554 ncdDetect we model the heterogeneous neutral background mutation-rate, taking genomic

555 annotations known to correlate with the mutation rate into account. We consider the mutational
556 burden and functional impact to reveal signs of recurrent positive selection across cancer
557 genomes.

558

559 The position- and sample-specific approach behind ncdDetect sets the stage for a number of
560 distinct types of analyses. The analysis of one contiguous region is a straight-forward
561 application of ncdDetect, as is the combined analysis of disjoint regions, potentially with vastly
562 different background mutation-rates. The flexible setup conveniently ensures that no constraints
563 are necessary when defining the size and location of a particular region of interest. Furthermore,
564 the method can be used to evaluate more complex functional hypotheses than those presented
565 here. For instance, different sets of regions in different samples can be jointly evaluated, and
566 sample- or tissue-specific scoring schemes can be applied directly.

567

568 Not all of the significant non-coding elements can be regarded as true cancer drivers. ncdDetect
569 might falsely identify driver elements (“false positives”) for both technical and biological reasons.
570 The false positives stemming from predicting too low a mutation rate in certain regions can be
571 reduced by adding relevant genomic annotations as explanatory variables to the null model. In
572 general, failing to include an explanatory variable that explains variation in the mutation rate will
573 cause too little variation in the predicted mutational probabilities. We handle such
574 overdispersion of the mutation rate by adjusting each sample- and position-specific probability
575 of mutation with an overdispersion-based correction factor. This improves the model fit,
576 although some inflation appears to remain for long elements. We acknowledge that the false
577 positive rate among long genes is not properly controlled and likely higher than the applied FDR
578 threshold. We therefore continue to strive to improve our model of the site-specific mutational
579 process, which will also improve power. The observed mutation rate also varies for technical
580 reasons. For instance, the power to call mutations and the rate with which mutations are

581 missed, will vary with genomic complexity, including repeats, pseudogenes, etc. The predicted
582 mutational probabilities can thus be further improved by including genomic annotations that
583 correlate with the rate of either false negative or false positive mutation calls as explanatory
584 variables.

585

586 Likewise, ncdDetect might miss true driver elements (“false negatives”). Especially, with the size
587 of the current whole cancer-genome data sets, we lack statistical power to detect infrequently
588 mutated driver elements, or driver elements that may operate within an individual cancer type.
589 This issue will be remedied as larger sets of sequenced whole genomes become available. In
590 the near future, more than 2,500 whole genomes will be available from the Pan-Cancer Analysis
591 of Whole Genomes (PCAWG) project. However, it is becoming evident that some instances of
592 detected potential driver regions may be explained by local mutational mechanisms rather than
593 recurrent selection (Sabarinathan et al. 2016). This emphasizes the importance of critical
594 scrutinization and eventually independent validation of driver candidates emerging from
595 ncdDetect.

596

597 With ncdDetect we screen for non-coding cancer drivers and highlight cases of special interest.
598 To gain further evidence for the identified candidates, we correlate the presence of non-coding
599 mutations with gene expression as well as patient survival: We find that mutations in the
600 promoter and in the coding region of a gene in the Base Excision Repair pathway, *SMUG1*,
601 correlate with an increase of C→T mutations. We hypothesise that *SMUG1* mutations might
602 perturb uracil-DNA glycosylase function and cause a specific mutational phenotype. Though our
603 study is limited to correlative observations between the expected mutational signature for uracil-
604 DNA glycosylase deficiency and mutational presence as well as gene expression, perturbation
605 experiments in cell lines support the hypothesis (An et al. 2005). We also identify non-coding

606 regulatory regions that associate with patient survival, including the potential clinically important
607 5' UTR of *CD1A*, the promoter and 3'UTR of *PRSS3*, and the 3'UTR of *SEC14L1*. Finally, we
608 identify lung cancer mutations in the splice sites of *STK11* as potential driver events. By
609 extending the analysis to the larger TCGA data set, we show that these correlate significantly
610 with expression. The patients also show a strong tendency for poorer survival.

611
612 In this work, we have addressed the challenges associated with distinguishing driver and
613 passenger non-coding mutations. We evaluated three different scoring schemes and found that
614 a conservation-based scheme performed better than mutation counts and log-likelihoods in our
615 setting. For selected candidate cases, we found a significant effect on expression levels and a
616 significant decrease in survival for mutated samples. The combined analyses of mutational
617 impact on expression and survival across cancer types allowed us aggregate evidence and gain
618 power. The screen identified candidates of potential clinical relevance. However, sample sizes
619 remain small and further studies in large independent cohorts are necessary to establish their
620 potential as prognostic biomarkers or therapeutic targets.

621
622 As we continue to gain larger cancer genomics data sets for driver screens, accurate modelling
623 of the mutational heterogeneity will become increasingly important. This will help control the
624 false positive rate as the power of the data increases. As our understanding of the general
625 differences in mutational mechanisms between cancer types improves further, this knowledge
626 should be incorporated in ncdDetect.

627 Materials and methods

628 Statistical null model

629 The statistical null model that enables us to predict position- and sample-specific probabilities of
630 mutation is a multinomial logistic regression model (Agresti 2013). The model is described in
631 detail in (Bertl et al., n.d.). Logistic regression has been used to model the background mutation
632 rate in cancer in previous non-coding driver detection studies (Melton et al. 2015). The model
633 considers four possible outcomes; transitions ($TS_{\{A \rightarrow G, G \rightarrow A\}}$), two types of transversions ($TV_{\{A \rightarrow T, G \rightarrow T\}}$ and $TV_{\{A \rightarrow C, G \rightarrow C\}}$) as well as the reference class of no mutation. The use of logistic
634 regression ensures that the predicted probabilities are restricted to lie in the interval between
635 zero and one. The reference sequence used in the model is the GRCh37 assembly (hg19) for
636 the human genome. The explanatory variables of the model are listed below.

638

- 639 • Sample id: a factor variable with 505 levels.
- 640 • Replication timing: A numeric variable with values between zero (early replication) and
641 one (late replication) (Chen et al. 2010). The variable is computed for 100kb windows.
642 Originally, the variable corresponds to the hg18 assembly for the human genome. It is
643 converted to the hg19 assembly using the UCSC liftOver tool (“UCSC Genome Browser”
644 2016).
- 645 • Trinucleotides: This variable is broken down into two separate variables; the reference
646 bp in question as well as the left and right flanking bases. The bp in question is encoded
647 as a factor variable with two levels, “A” for A:T bps with weak hydrogen bonds, and “G”
648 for G:C bps with strong hydrogen bonds. The left and right flanking bases are

649 implemented in a factor variable with 16 different levels, “AA” to “TT”. Including this
650 variable in an interaction term with sample id effectively takes the sample-specific
651 mutational signatures into account.

- 652 • Genomic segment: A factor variable with six levels, “protein-coding genes”, “promoter
653 elements”, “splice sites”, “5’ UTRs”, “3’ UTRs” and “other”.
- 654 • Expression level: A numeric variable based on all available RNAseq expression data
655 from TCGA (version 2, RSEM values, level 3 data). All RSEM values were $\log_2(x+1)$
656 transformed. For each cancer type, the median expression was calculated for all genes.
657 If multiple annotations of a gene existed, the longest annotation was used. For
658 overlapping genes, the expression is summed up. We collapsed colon (COAD) and
659 rectal carcinoma (READ) to a joint cancer type CRC by averaging over the expression
660 values (Fredriksson et al. 2014).
- 661 • Local mutation rate: A numeric variable calculated per base position. For each position,
662 the position itself plus the flanking 10kb on either side is skipped to avoid that the rate of
663 the tested element has a large effect on local mutation rate. The local mutation rate is
664 then based on the next flanking 20kb regions on either side of the skipped regions. For
665 each sample, the number of mutations in the two 20kb regions are weighted by the
666 reciprocal of the total number of mutations in the sample. The value of the local mutation
667 rate is the weighted sum of the mutations across all samples.

668

669 The multinomial logistic regression model fit is based on a so-called *count-table*. For the
670 purpose of creating this data structure, the three numeric variables, replication timing, local
671 mutation rate, and expression level are each discretized into five bins. For each combination of
672 explanatory variable levels ($505 \times 5 \times 2 \times 16 \times 6 \times 5 \times 5 = 12,120,000$ combinations), the

673 number of genomic positions as well as the number of mutations of each type are counted.
674 Before constructing the count-table, all COSMIC genes were excluded from the data set.
675
676 For fast and memory-efficient estimation, the multinomial logistic regression model is split up
677 into three binomial logistic models (Begg and Gray 1984). Estimation is conducted in R (R
678 Development Core Team 2008) (RRID:SCR_001905) using the function *glm4* from the
679 contributed package *MatrixModels* (Bates and Maechler 2015), which provides efficient
680 estimation for GLMs with sparse design matrices. Three-fold multiple imputation is used to
681 handle missing values in the variable replication timing (Schafer 1997). The imputed values are
682 randomly drawn from the marginal distribution of observed replication timing values.

683 Scoring schemes

684 In the implementation of ncdDetect, the scores must be discrete values. For speed efficiency,
685 integer values are recommended. The conservation scoring scheme is based on the phyloP
686 scores. To ensure that all scores are positive, the phyloP values are shifted by adding 20 to all
687 values. For computational reasons, the values are rounded downwards to the first decimal point,
688 and multiplied by a factor of 10 to create integers. No mutation is associated with a score value
689 of zero. The log-likelihood scoring scheme defines the scores as minus the natural logarithm of
690 the sample- and position-specific neutral somatic mutation probabilities predicted by the null
691 model. The scores are converted into integer values by the same procedure used for the
692 conservation scores. Effectively, this means that positions with no mutations will be given a
693 score of zero.

694 Candidate elements

695 The candidate elements are defined based on the protein-coding transcript annotations of
696 GENCODE version 19 basic annotation set (GENCODE 2016). Regions are divided into five
697 categories; protein-coding genes, promoter elements, splice sites, 3' UTRs and 5' UTRs. The
698 different regions are defined per transcript and collapsed per gene. Splice site regions are
699 defined as the two intronic bases on either side of all internal exons. Promoter elements are
700 defined as 500 bps in either direction from transcriptional start sites. A hierarchy for the
701 categories is defined as protein-coding genes > splice sites > 3' UTRs > 5' UTRs > promoter
702 elements. Bps included in two or more categories are retained only for the category higher in
703 the hierarchy. Region definitions are available in Supplementary File 4.

704

705 Elements located on chromosome X and Y are not considered in the analyses conducted here.
706 The number of analysed elements thus differ from the total number of elements defined (Table
707 1). We analyse 19,256 protein-coding regions, 19,157 promoter elements, 17,867 splice sites,
708 18,481 3' UTRs and 18,220 5' UTRs. The length distributions of the element types are depicted
709 in Figure 5-figure supplement 3.

710 An overdispersion-based mutation rate adjustment

711 To correct for overdispersion of the mutation rate, we adjust each sample- and position-specific
712 mutational probability by an overdispersion-based mutation rate correction factor. The correction
713 factor is modelled with a beta binomial model. Details are given in Appendix, section 1.

714 TCGA exome data

715 To support our findings from the 505 whole genome TCGA samples, we obtain mutation calls,
716 expression values and survival data based on the larger set of TCGA exomes (Cancer Genome

717 Atlas Research Network et al. 2013, Supplementary File 7). These data are applied in the
718 expression analyses, in the analyses of enrichment of G→A mutations in *SMUG1* mutated
719 samples, as well as in the survival analyses performed.

720

721 We obtain mutation calls for 5,802 samples. Of these, we remove 348 samples, which are also
722 present in the original 505-sample data set. The final mutation call set thus consists of 5,454
723 TCGA exome samples. We obtain expression data for 8,471 samples, and after removing
724 samples present in the 505-sample set, a total of 4,295 samples have both mutation calls and
725 expression data available. A total of 5,336 TCGA exome samples have both mutation calls and
726 clinical survival data available, after subtracting samples that are also present in the 505 whole
727 genome TCGA sample set.

728 Expression analysis

729 For a given candidate, we perform a two-sided Wilcoxon rank sum test. With this, we test the
730 hypothesis that there is no difference in gene expression levels between samples that are
731 mutated and samples that are not mutated in a given candidate element. Such a test is
732 performed for each individual cancer type, and the p-values are combined across cancer types
733 using Fisher's method.

734

735 These analyses are performed for both the 505 whole genome TCGA samples and the 4,295
736 TCGA exome samples with the necessary data available. Data overview and results are
737 available in Supplementary File 5.

738 Enrichment of G→A mutations in *SMUG1* mutated samples

739 In order to test if the proportion of G→A (including C→T) mutations outside CpG sites are
740 greater for samples harboring a *SMUG1* mutation, compared to samples not carrying such a
741 mutation, we first count the 192 (= $4 \times 4 \times 4 \times 3$) different mutation types (including left and
742 right flanking bases for a given mutated position) for each of the samples. For each sample, we
743 count the number of G→A mutations, excluding those in CpG sites and, for melanoma samples,
744 those part of the CC→TT UV induced mutational signature (Alexandrov and Stratton 2014). The
745 counts are normalized by the total number of mutations for the sample. These proportions are
746 referred to as the *uracil-DNA glycosylase deficiency signature*. For a given cancer type, we
747 divide the samples into two groups; those that carry a *SMUG1* mutation, and those that do not.
748 A one-sided Wilcoxon rank sum test is performed to test the null hypothesis that *SMUG1*
749 mutated samples do not have higher values of the uracil-DNA glycosylase deficiency signature
750 statistic, compared to samples without *SMUG1* mutations. Fisher's method is used to combine
751 p-values across cancer types. This type of analysis is performed for the 505 whole genome
752 TCGA sample set.

753

754 We further analyse whether there is a correlation between the uracil-DNA glycosylase
755 deficiency signature and *SMUG1*, *UNG*, or *SMUG1* × *UNG* gene expression. This type of
756 analysis is performed for both the 505 whole genome TCGA samples and the 4,295 TCGA
757 exome samples with the necessary data available. The uracil-DNA glycosylase deficiency

758 signature statistic calculated on the basis of the TCGA exome data set is based only on
759 captured CDS regions. Results and data overview are available in Supplementary File 5.

760 Survival analysis

761 We investigate the correlation between mutation status and survival data in the following
762 manner: We download clinical data from the TCGA data portal ("TCGA Data Portal" 2016,
763 Supplementary File 7) (running date 01/11/2015) using the RTCGAToolbox R library (Samur
764 2014). For a candidate driver element, the difference in survival between mutated and non-
765 mutated samples is tested per cancer type using a one-sided Log-rank test on the Kaplan-Meier
766 estimated survival curves (Kaplan and Meier 1958). The tests are only performed when at least
767 four mutations are observed within a given cancer type. Evidence is combined across cancer
768 types with Fisher's method.

769
770 The survival analysis is performed for each of the top-50 ranked ncdDetect candidates of each
771 non-coding element type (promoters, splice sites, 3' UTRs and 5' UTRs), or all significant
772 elements of a given type. The analyses are performed for both the 505 whole genome TCGA
773 sample set, and the 5,336 TCGA exome sample set with mutation calls and clinical survival data
774 available. Results and data overview are available in Supplementary File 6.

775 Time complexity

776 The use of mathematical convolution on the fine grained sample- and position-specific scores
777 and probabilities makes ncdDetect computationally intensive (Figure 3-figure supplement 1).
778 Convolution is the procedure of calculating the distribution function of the sum of independent
779 discrete random variables. The algorithm is implemented using dynamic programming (Touzet
780 and Varré 2007), and can be thought of as filling out a matrix from the bottom left corner to the

781 upper right corner. Convoluting each cell in the matrix has time complexity $O(1)$, and the running
782 time is thus determined by the size of the matrix. The time complexity of the algorithm is $O(m \times$
783 $k \times s_{\max})$, where m is the element size, k is the number of samples and s_{\max} is the maximum
784 score.

785 Implementation

786 `ncdDetect` is implemented in the software environment R (R Development Core Team 2008)
787 (RRID:SCR_001905), using the *Rcpp* and *RcppArmadillo* packages (Eddelbuettel 2016) for
788 speed optimization. The core `ncdDetect` functions to perform convolution are collected in the R
789 package *ncdDetectTools* available at github.com. The package can be installed using the
790 *devtools* package (“Tools to Make Developing R Packages Easier [R Package Devtools Version
791 1.12.0]” 2016): `install_github("MaleneJuul/ncdDetectTools")`. A few examples of the package
792 functionalities are given in the package vignette, also available in the github repository.
793 The null model estimates used for the current application is provided at
794 <http://moma.ki.au.dk/ncddetect/> along with a tutorial on how to obtain p-values from these
795 estimates using `ncdDetect`.

796 Availability of data and materials

797 Mutational, expression and clinical data for the TCGA samples are administered by dbGaP
798 (<https://dbgap.ncbi.nlm.nih.gov>) (RRID:SCR_002709). The additional datasets supporting the
799 conclusions of this article are included within the article and its supplementary files.

800 Acknowledgements

801 We thank the TCGA consortium for data access and the system administrators of the
802 GenomeDK high performance computing facility for facilitating the computational analysis.

803 Competing interests

804 The authors declare that they have no competing interests.

805 References

806 Agresti, Alan. 2013. *Categorical Data Analysis*. John Wiley & Sons.
807 Alexandrov, Ludmil B., Nik-Zainal Serena, David C. Wedge, Samuel A J, Behjati Sam, Andrew
808 V. Biankin, Graham R. Bignell, et al. 2013. "Signatures of Mutational Processes in Human
809 Cancer." *Nature* 500 (7463): 415–21.
810 Alexandrov, Ludmil B., and Michael R. Stratton. 2014. "Mutational Signatures: The Patterns of
811 Somatic Mutations Hidden in Cancer Genomes." *Current Opinion in Genetics &
812 Development* 24 (February): 52–60.
813 An, Qian, Peter Robins, Tomas Lindahl, and Deborah E. Barnes. 2005. "C → T Mutagenesis
814 and γ -Radiation Sensitivity due to Deficiency in the Smug1 and Ung DNA Glycosylases."
815 *The EMBO Journal* 24 (12): 2205–13.
816 Bates, Douglas, and Martin Maechler. 2015. "MatrixModels: Modelling with Sparse And Dense
817 Matrices." <http://CRAN.R-project.org/package=MatrixModels>.
818 Begg, Colin B., and Robert Gray. 1984. "Calculation of Polychotomous Logistic Regression
819 Parameters Using Individualized Regressions." *Biometrika* 71 (1): 11–18.
820 Benjamini, Yoav, and Yosef Hochberg. 1995. "Controlling the False Discovery Rate: A Practical
821 and Powerful Approach to Multiple Testing." *Journal of the Royal Statistical Society. Series
822 B, Statistical Methodology* 57 (1). [Royal Statistical Society, Wiley]: 289–300.
823 Bertl, Johanna, Qianyun Guo, Malene Juul, Søren Besenbacher, Asger Hobolth, and Jakob
824 Skou Pedersen. n.d. "A Site Specific Model and Analysis of the Neutral Somatic Mutation
825 Rate in Whole-Genome Cancer Data." *In Preparation (expected Publication at Biorxiv.org
826 Mid March 2017)*.
827 Burdelski, Christoph, Barreau Ysé, Simon Ronald, Hube-Magg Claudia, Minner Sarah, Koop
828 Christina, et al. 2015. "Saccharomyces Cerevisiae-like 1 Overexpression Is Frequent in
829 Prostate Cancer and Has Markedly Different Effects in Ets-Related Gene Fusion-positive
830 and Fusion-Negative Cancers." *Human Pathology* 46 (4): 514–23.
831 Cancer Genome Atlas Research Network, John N. Weinstein, Eric A. Collisson, Gordon B. Mills,
832 Kenna R. Mills Shaw, Brad A. Ozenberger, Kyle Ellrott, Ilya Shmulevich, Chris Sander, and
833 Joshua M. Stuart. 2013. "The Cancer Genome Atlas Pan-Cancer Analysis Project." *Nature
834 Genetics* 45 (10): 1113–20.
835 Chen, Chun-Long, Aurélien Rappailles, Lauranne Duquenne, Maxime Huvet, Guillaume
836 Guilbaud, Laurent Farinelli, Benjamin Audit, et al. 2010. "Impact of Replication Timing on
837 Non-CpG and CpG Substitution Rates in Mammalian Genomes." *Genome Research* 20 (4):
838 447–57.
839 Cornish-Bowden, Athel. 1985. "Nomenclature for Incompletely Specified Bases in Nucleic Acid
840 Sequences: Recommendations 1984." *Nucleic Acids Research* 13 (9): 3021–30.
841 Cosmic. 2016. "COSMIC: Cancer Gene Census." Accessed May 26.
842 <http://cancer.sanger.ac.uk/census/>.
843 Coventry, Brendon, and Susanne Heinzl. 2004. "CD1a in Human Cancers: A New Role for an
844 Old Molecule." *Trends in Immunology* 25 (5): 242–48.
845 Dees, N. D., Q. Zhang, C. Kandoth, M. C. Wendl, W. Schierding, D. C. Koboldt, T. B. Mooney,
846 et al. 2012. "MuSiC: Identifying Mutational Significance in Cancer Genomes." *Genome
847 Research* 22 (8): 1589–98.
848 Diederichs, Sven, Lorenz Bartsch, Julia C. Berkmann, Karin Fröse, Jana Heitmann, Caroline
849 Hoppe, Deetje Iggena, et al. 2016. "The Dark Matter of the Cancer Genome: Aberrations in
850 Regulatory Elements, Untranslated Regions, Splice Sites, Non-Coding RNA and
851 Synonymous Mutations." *EMBO Molecular Medicine* 8 (5): 442–57.

852 Ding, Jiarui, Melissa K. McConechy, Hugo M. Horlings, Gavin Ha, Fong Chun Chan, Tyler
853 Funnell, Sarah C. Mullaly, et al. 2015. "Systematic Analysis of Somatic Mutations Impacting
854 Gene Expression in 12 Tumour Types." *Nature Communications* 6 (October): 8554.
855 Drake, Charles G., Evan J. Lipson, and Julie R. Brahmer. 2014. "Breathing New Life into
856 Immunotherapy: Review of Melanoma, Lung and Kidney Cancer." *Nature Reviews. Clinical
857 Oncology* 11 (1): 24–37.
858 Edelbuettel, Dirk. 2016. "Rcpp: Seamless R and C++ Integration." Accessed June 14.
859 <http://www.rcpp.org>.
860 Ehrenberg, Philip K., Aviva Geretz, Karen M. Baldwin, Richard Apps, Victoria R. Polonis, Merlin
861 L. Robb, Jerome H. Kim, Nelson L. Michael, and Rasmi Thomas. 2014. "High-Throughput
862 Multiplex HLA Genotyping by next-Generation Sequencing Using Multi-Locus Individual
863 Tagging." *BMC Genomics* 15 (October): 864.
864 Eicheler, Wolfgang, Daniel Zips, Annegret Dörfler, Reidar Grénman, and Michael Baumann.
865 2002. "Splicing Mutations in TP53 in Human Squamous Cell Carcinoma Lines Influence
866 Immunohistochemical Detection." *The Journal of Histochemistry and Cytochemistry: Official
867 Journal of the Histochemistry Society* 50 (2): 197–204.
868 ENCODE Project Consortium. 2012. "An Integrated Encyclopedia of DNA Elements in the
869 Human Genome." *Nature* 489 (7414): 57–74.
870 Forbes, S. A., D. Beare, P. Gunasekaran, K. Leung, N. Bindal, H. Boutselakis, M. Ding, et al.
871 2014. "COSMIC: Exploring the World's Knowledge of Somatic Mutations in Human
872 Cancer." *Nucleic Acids Research* 43 (D1): D805–11.
873 Fredriksson, Nils J., Lars Ny, Jonas A. Nilsson, and Erik Larsson. 2014. "Systematic Analysis of
874 Noncoding Somatic Mutations and Gene Expression Alterations across 14 Tumor Types."
875 *Nature Genetics* 46 (12): 1258–63.
876 GENCODE. 2016. "GENCODE - Release History." Accessed September 8.
877 <http://www.encodegenes.org/releases/>.
878 Gerstein, Mark B., Anshul Kundaje, Manoj Hariharan, Stephen G. Landt, Koon-Kiu Yan, Chao
879 Cheng, Xinmeng Jasmine Mu, et al. 2012. "Architecture of the Human Regulatory Network
880 Derived from ENCODE Data." *Nature* 489 (7414): 91–100.
881 Gill, R. K., S-H Yang, D. Meerzaman, L. E. Mechanic, E. D. Bowman, H-S Jeon, S. Roy
882 Chowdhuri, et al. 2011. "Frequent Homozygous Deletion of the LKB1/STK11 Gene in Non-
883 Small Cell Lung Cancer." *Oncogene* 30 (35): 3784–91.
884 Gonzalez-Perez, Abel, and Nuria Lopez-Bigas. 2012. "Functional Impact Bias Reveals Cancer
885 Drivers." *Nucleic Acids Research* 40 (21): e169.
886 Grinstead, Charles Miller, and James Laurie Snell. 1997. *Introduction to Probability*. American
887 Mathematical Soc.
888 GTEx Consortium. 2013. "The Genotype-Tissue Expression (GTEx) Project." *Nature Genetics*
889 45 (6): 580–85.
890 Heidenreich, Barbara, Heidenreich Barbara, P. Sivaramakrishna Rachakonda, Hemminki Kari,
891 and Kumar Rajiv. 2014. "TERT Promoter Mutations in Cancer Development." *Current
892 Opinion in Genetics & Development* 24: 30–37.
893 Horn, Susanne, Adina Figl, P. Sivaramakrishna Rachakonda, Christine Fischer, Antje Sucker,
894 Andreas Gast, Stephanie Kadel, et al. 2013. "TERT Promoter Mutations in Familial and
895 Sporadic Melanoma." *Science* 339 (6122): 959–61.
896 Huang, Franklin W., Eran Hodis, Mary Jue Xu, Gregory V. Kryukov, Lynda Chin, and Levi A.
897 Garraway. 2013. "Highly Recurrent TERT Promoter Mutations in Human Melanoma."
898 *Science* 339 (6122): 957–59.
899 International Cancer Genome Consortium, Thomas J. Hudson, Warwick Anderson, Axel Artez,
900 Anna D. Barker, Cindy Bell, Rosa R. Bernabé, et al. 2010. "International Network of Cancer
901 Genome Projects." *Nature* 464 (7291): 993–98.
902 Jiang, Guozhong, Fengyu Cao, Guoping Ren, Dongling Gao, Vipul Bhakta, Yunhan Zhang, Hua

903 Cao, et al. 2010. "PRSS3 Promotes Tumour Growth and Metastasis of Human Pancreatic
904 Cancer." *Gut* 59 (11): 1535–44.

905 Kalikin, L. M., E. M. Bugeaud, P. L. Palmboos, R. H. Lyons Jr, and E. M. Petty. 2001. "Genomic
906 Characterization of Human SEC14L1 Splice Variants within a 17q25 Candidate Tumor
907 Suppressor Gene Region and Identification of an Unrelated Embedded Expressed
908 Sequence Tag." *Mammalian Genome: Official Journal of the International Mammalian
909 Genome Society* 12 (12): 925–29.

910 Kaplan, E. L., and Paul Meier. 1958. "Nonparametric Estimation from Incomplete Observations."
911 *Journal of the American Statistical Association* 53 (282): 457.

912 Kent, W. J. 2002. "The Human Genome Browser at UCSC." *Genome Research* 12 (6): 996–
913 1006.

914 Khurana, Ekta, Yao Fu, Dimple Chakravarty, Francesca Demichelis, Mark A. Rubin, and Mark
915 Gerstein. 2016. "Role of Non-Coding Sequence Variants in Cancer." *Nature Reviews.
916 Genetics* 17 (2): 93–108.

917 Khurana, Ekta, Yao Fu, Vincenza Colonna, Xinmeng Jasmine Mu, Hyun Min Kang, Tuuli
918 Lappalainen, Andrea Sboner, et al. 2013. "Integrative Annotation of Variants from 1092
919 Humans: Application to Cancer Genomics." *Science* 342 (6154): 1235587.

920 Lanzós, Andrés, Joana Carlevaro-Fita, Loris Mularoni, Ferran Reverter, Emilio Palumbo,
921 Roderic Guigó, and Rory Johnson. 2017. "Discovery of Cancer Driver Long Noncoding
922 RNAs across 1112 Tumour Genomes: New Candidates and Distinguishing Features."
923 *Scientific Reports* 7 (January): 41544.

924 Lawrence, Michael S., Petar Stojanov, Paz Polak, Gregory V. Kryukov, Kristian Cibulskis,
925 Andrey Sivachenko, Scott L. Carter, et al. 2013. "Mutational Heterogeneity in Cancer and
926 the Search for New Cancer-Associated Genes." *Nature* 499 (7457): 214–18.

927 Lee, Eung Bae, Guang Jin, Shin Yup Lee, Ji Young Park, Min Jung Kim, Jin Eun Choi, Hyo
928 Sung Jeon, et al. 2010. "TP53 Mutations in Korean Patients with Non-Small Cell Lung
929 Cancer." *Journal of Korean Medical Science* 25 (5): 698–705.

930 Lochovsky, Lucas, Jing Zhang, Yao Fu, Ekta Khurana, and Mark Gerstein. 2015. "LARVA: An
931 Integrative Framework for Large-Scale Analysis of Recurrent Variants in Noncoding
932 Annotations." *Nucleic Acids Research* 43 (17): 8123–34.

933 Ma, Ruiqiong, Xue Ye, Hongyan Cheng, Yu Ma, Heng Cui, and Xiaohong Chang. 2015.
934 "PRSS3 Expression Is Associated with Tumor Progression and Poor Prognosis in Epithelial
935 Ovarian Cancer." *Gynecologic Oncology* 137 (3): 546–52.

936 Marx, Vivien. 2014. "Cancer Genomes: Discerning Drivers from Passengers." *Nature Methods*
937 11 (4): 375–79.

938 Melton, Collin, Jason A. Reuter, Damek V. Spacek, and Michael Snyder. 2015. "Recurrent
939 Somatic Mutations in Regulatory Regions of Human Cancer Genomes." *Nature Genetics*
940 47 (7): 710–16.

941 Mularoni, Loris, Radhakrishnan Sabarinathan, Jordi Deu-Pons, Abel Gonzalez-Perez, and Núria
942 López-Bigas. 2016. "OncodriveFML: A General Framework to Identify Coding and Non-
943 Coding Regions with Cancer Driver Mutations." *Genome Biology* 17 (1): 128.

944 Pettersen, Henrik Sahlin, Galashevskaya Anastasia, Doseth Berit, Mirta M. L. Sousa, Sarno
945 Antonio, Visnes Torkild, Per Arne Aas, et al. 2015. "AID Expression in B-Cell Lymphomas
946 Causes Accumulation of Genomic Uracil and a Distinct AID Mutational Signature." *DNA
947 Repair* 25: 60–71.

948 Pettersen, Henrik Sahlin, Ottar Sundheim, Karin Margaretha Gilljam, Geir Slupphaug, Hans
949 Einar Krokan, and Bodil Kavli. 2007. "Uracil-DNA Glycosylases SMUG1 and UNG2
950 Coordinate the Initial Steps of Base Excision Repair by Distinct Mechanisms." *Nucleic
951 Acids Research* 35 (12): 3879–92.

952 Polak, Paz, Rosa Karlić, Amnon Koren, Robert Thurman, Richard Sandstrom, Michael S.
953 Lawrence, Alex Reynolds, et al. 2015. "Cell-of-Origin Chromatin Organization Shapes the

954 Mutational Landscape of Cancer.” *Nature* 518 (7539): 360–64.

955 Pollard, Katherine S., Melissa J. Hubisz, Kate R. Rosenbloom, and Adam Siepel. 2010.

956 “Detection of Nonneutral Substitution Rates on Mammalian Phylogenies.” *Genome*

957 *Research* 20 (1): 110–21.

958 Pon, Julia R., and Marco A. Marra. 2015. “Driver and Passenger Mutations in Cancer.” *Annual*

959 *Review of Pathology: Mechanisms of Disease* 10 (1): 25–50.

960 R Development Core Team. 2008. “R: A Language and Environment for Statistical Computing.”

961 Vienna, Austria: R Foundation for Statistical Computing. <http://www.R-project.org>.

962 Reimand, Jüri, and Gary D. Bader. 2013. “Systematic Analysis of Somatic Mutations in

963 Phosphorylation Signaling Predicts Novel Cancer Drivers.” *Molecular Systems Biology* 9:

964 637.

965 Reva, Boris, Yevgeniy Antipin, and Chris Sander. 2011. “Predicting the Functional Impact of

966 Protein Mutations: Application to Cancer Genomics.” *Nucleic Acids Research* 39 (17):

967 e118.

968 Sabarinathan, Radhakrishnan, Loris Mularoni, Jordi Deu-Pons, Abel Gonzalez-Perez, and Núria

969 López-Bigas. 2016. “Nucleotide Excision Repair Is Impaired by Binding of Transcription

970 Factors to DNA.” *Nature* 532 (7598): 264–67.

971 Samur, Mehmet Kemal. 2014. “RTCGAToolbox: A New Tool for Exporting TCGA Firehose

972 Data.” *PloS One* 9 (9): e106397.

973 Schafer, J. L. 1997. *Analysis of Incomplete Multivariate Data*. CRC Press.

974 Schweder, T., and E. Spjøtvoll. 1982. “Plots of P-Values to Evaluate Many Tests

975 Simultaneously.” *Biometrika* 69 (3): 493.

976 Srebrow, Anabella, and Alberto R. Kornblihtt. 2006. “The Connection between Splicing and

977 Cancer.” *Journal of Cell Science* 119 (Pt 13): 2635–41.

978 Stratton, Michael R., Peter J. Campbell, and P. Andrew Futreal. 2009. “The Cancer Genome.”

979 *Nature* 458 (7239): 719–24.

980 Tamborero, David, Abel Gonzalez-Perez, and Nuria Lopez-Bigas. 2013. “OncodriveCLUST:

981 Exploiting the Positional Clustering of Somatic Mutations to Identify Cancer Genes.”

982 *Bioinformatics* 29 (18): 2238–44.

983 “TCGA Data Portal.” 2016. Accessed July 15. [https://ocg.cancer.gov/resources/cancer-genome-](https://ocg.cancer.gov/resources/cancer-genome-atlas-tcga-data-portal)

984 [atlas-tcga-data-portal](https://ocg.cancer.gov/resources/cancer-genome-atlas-tcga-data-portal).

985 “Tools to Make Developing R Packages Easier [R Package Devtools Version 1.12.0].” 2016.

986 Comprehensive R Archive Network (CRAN). Accessed August 24. [http://CRAN.R-](http://CRAN.R-project.org/package=devtools)

987 [project.org/package=devtools](http://CRAN.R-project.org/package=devtools).

988 Touzet, H el ene, and Jean-St ephane Varr e. 2007. “Efficient and Accurate P-Value Computation

989 for Position Weight Matrices.” *Algorithms for Molecular Biology: AMB* 2 (December): 15.

990 “UCSC Genome Browser.” 2016. Accessed May 20.

991 <http://hgdownload.cse.ucsc.edu/downloads.html>.

992 Van Rhijn, Ildiko, Dale I. Godfrey, Jamie Rossjohn, and D. Branch Moody. 2015. “Lipid and

993 Small-Molecule Display by CD1 and MR1.” *Nature Reviews. Immunology* 15 (10): 643–54.

994 Varley, J. M., C. Attwooll, G. White, G. McGown, M. Thorncroft, A. M. Kelsey, M. Greaves, J.

995 Boyle, and J. M. Birch. 2001. “Characterization of Germline TP53 Splicing Mutations and

996 Their Genetic and Functional Analysis.” *Oncogene* 20 (21): 2647–54.

997 Vinagre, Jo ao, Vinagre Jo ao, Almeida Ana, P opulo Helena, Batista Rui, Lyra Joana, Pinto

998 Vasco, et al. 2013. “Frequency of TERT Promoter Mutations in Human Cancers.” *Nature*

999 *Communications* 4. doi:10.1038/ncomms3185.

1000 Visnes, Torkild, Berit Doseth, Henrik Sahlin Pettersen, Lars Hagen, Mirta M. L. Sousa, Mansour

1001 Akbari, Marit Otterlei, Bodil Kavli, Geir Slupphaug, and Hans E. Krokan. 2009. “Uracil in

1002 DNA and Its Processing by Different DNA Glycosylases.” *Philosophical Transactions of the*

1003 *Royal Society of London. Series B, Biological Sciences* 364 (1517): 563–68.

1004 Weinhold, Nils, Anders Jacobsen, Nikolaus Schultz, Chris Sander, and William Lee. 2014.

1005 "Genome-Wide Analysis of Noncoding Regulatory Mutations in Cancer." *Nature Genetics*
1006 46 (11): 1160–65.
1007 Wei, Shuanzeng, Virginia A. LiVolsi, Marcia S. Brose, Kathleen T. Montone, Jennifer J. D.
1008 Morrissette, and Zubair W. Baloch. 2016. "STK11 Mutation Identified in Thyroid
1009 Carcinoma." *Endocrine Pathology* 27 (1): 65–69.
1010 Yang, P., C. W. Du, M. Kwan, S. X. Liang, and G. J. Zhang. 2013. "The Impact of p53 in
1011 Predicting Clinical Outcome of Breast Cancer Patients with Visceral Metastasis." *Scientific*
1012 *Reports* 3: 2246.

1013

1014 Figure legends

1015 Figure 1

1016 Variation in mutation rate at different scales and various explanatory variables. **A:** The flowchart
1017 illustrates the input to the model fit that predicts the position- and sample-specific mutational
1018 probabilities. **B:** The number of mutations observed per sample divided into the 14 different
1019 cancer types. **C:** The set of genomic annotations used as explanatory variables are illustrated
1020 on a 300 kb region of chromosome 1 for the colorectal cancer sample CRC_TCGA-A6-6141-
1021 01A. For illustrative purposes, the nucleotide sequence is shown on a 30-bp section of
1022 chromosome 1 and trinucleotides likewise on a 5-bp section.

1023 Figure 2

1024 Position- and sample-specific predicted mutation rates and scoring-schemes. **A:** A multinomial
1025 logistic regression model is used to predict the sample- and position-specific background
1026 mutation-probabilities. **B:** The genomic annotations and the reference sequence (Figure 1) are
1027 used as explanatory variables in a regression fit of the somatic mutation rate. In effect, a logistic
1028 regression model is fitted for each of the four types of outcome (three types of mutation and no
1029 mutation) and combined into a multinomial logistic regression fit. Logistic regression ensures
1030 probability-predictions between zero and one. The mutation probabilities are of such small

1031 magnitude that we observe near linearity of the logistic regression curve. **C:** Sample- and
1032 position-specific predicted mutation probabilities for each of the four outcomes in a 300-bp
1033 region of chromosome 1 (chr1:115,824,535-115,824,834) for the colorectal cancer sample
1034 CRC_TCGA-A6-6141-01A. **D:** Observed sample-specific somatic mutations within the same
1035 region. For the sample in question, two mutations are observed; one of type $TV_{\{A \rightarrow T, G \rightarrow T\}}$ and
1036 one of type $TV_{\{A \rightarrow C, G \rightarrow C\}}$. **E:** Sample- and position-specific scores for each of the three
1037 considered scoring schemes.

1038 Figure 3

1039 ncdDetect analysis concepts. **A:** Flowchart of the algorithmic steps of ncdDetect. Panels B
1040 through E show the sample-specific calculations, while panels F and G show the calculations
1041 across samples. **B:** The genomic candidate region is annotated with position- and sample-
1042 specific scores. The values of these scores depend on the choice of scoring scheme. **C:** The
1043 region is also annotated with sample- and position-specific predicted mutation probabilities.
1044 These probabilities are predicted by the null model, and does not depend on the choice of
1045 scoring scheme. **D:** The observed score of the sample is defined as the sum of the scores
1046 associated with the observed mutational events. Scores based on number of mutations and
1047 conservation will assign non-mutated positions with a score-value of zero. Scores based on log-
1048 likelihoods will assign non-mutated positions with a positive score-value, which in practice will
1049 be near zero. **E:** The sample-specific background score-distribution is obtained by convolution.
1050 **F:** Sample-specific calculations are carried out for each individual sample in the dataset. **G:** The
1051 overall background score-distribution is obtained by convolution of the individual-sample
1052 distributions. This figure is conceptual and not based on actual data. Figure 4D-F are real
1053 examples of background score-distributions.

1054 **Figure 4**

1055 Analysis of protein-coding genes to evaluate ncdDetect performance. **A:** The final null model is
1056 obtained through forward model-selection. The QQ-plot shows the p-values of all genes (n =
1057 19,256) plotted against their uniform expectation under the null for each of the five models
1058 considered. Deviations from the expectations (red identity line) are seen for a varying proportion
1059 of the genes (0.5-10%). Results are shown for conservation scores. Similar plots for log-
1060 likelihoods and number of mutations are shown in Figure 4-figure supplement 1. **B:** Venn
1061 diagram showing the overlap between protein-coding genes called as drivers by ncdDetect ($q <$
1062 0.10) for the three scoring schemes and the COSMIC Gene Census list. **C:** COSMIC Gene
1063 Census recall plot. The fraction of COSMIC genes recalled in the top ncdDetect candidates. **D-**
1064 **F:** The two most significant genes called by ncdDetect are *TP53* and *PIK3CA*. An example of a
1065 gene not called significant is *SLFN11*. For each of these, the convoluted background score-
1066 distributions are shown together with the observed scores and resulting p-values.

1067 **Figure 5**

1068 Q-values and top-ten ranking non-coding elements for each of the three proposed scoring
1069 schemes. The results discussed in the text relate to conservation scores. Non-coding elements
1070 associated to COSMIC genes are highlighted in red. For each element, the region size is given
1071 together with the observed number of mutations and the expected number of mutations under
1072 the null model. **A:** The QQ-plot shows the p-values for all promoter elements (n = 19,157)
1073 plotted against their uniform expectation under the null. 160 promoter elements are found to be
1074 significant. **B:** QQ-plot of p-values for all splice sites (n = 17,867). The p-values do not follow the
1075 expectation under the null. This is explained by the fact that 90% of all splice sites carry no
1076 mutations. Three splice sites come up significant with ncdDetect after correcting for multiple
1077 testing.

1078 **Figure 6**

1079 *SMUG1* mutations and base excision repair. **A:** Genomic overview of *SMUG1* showing its
1080 promoter region (Kent 2002). The DNase clusters track shows DNase hypersensitive regions
1081 where the darkness is proportional to the maximum signal strength observed in any cell line
1082 (ENCODE Project Consortium 2012). The transcription factor binding sites (TFBSs) track shows
1083 core regions of transcription factor binding (Gerstein et al. 2012). The phyloP track shows
1084 evolutionary conservation of positions (Pollard et al. 2010). **B:** Uracil-DNA glycosylase
1085 deficiency signature definition: (1) Cytosines may be methylated (orange circle) at CpG sites
1086 (gray box). (2) Spontaneous deamination (red boxes) of non-methylated cytosine results in
1087 uracil, causing U:G mismatches. Spontaneous deamination of methylated cytosine results in
1088 thymine, causing T:G mismatches. (3a) *SMUG1* and *UNG* are uracil-DNA glycosylases, which,
1089 via base excision repair, will repair the U:G mismatches caused by deamination. (3b) If
1090 unrepaired, the U:G mismatches will result in G→A mutations. **C:** A one-sided Wilcoxon rank
1091 sum test is performed per cancer type to investigate if samples with a *SMUG1* mutation have a
1092 higher value of the uracil-DNA glycosylase deficiency signature statistic than samples without
1093 such a mutation. The analysis is based on the 505 whole genome TCGA samples. Each dot
1094 represents a sample, and the color represents the *SMUG1*-associated mutated element. **D:**
1095 Correlation between the uracil-DNA glycosylase deficiency signature statistic and the product of
1096 *SMUG1* and *UNG* gene expression using TCGA exome data for lung adenocarcinoma.

1097 **Figure 7**

1098 Survival- and expression analysis of *CD1A*, *PRSS3* and *STK11* mutations. **A:** Kaplan-Meier
1099 survival curves for Melanoma samples with and without mutations in the 5' UTR of *CD1A*. For
1100 illustration purposes, the data is shown for a follow-up time of 2,000 days, at which point 98 out

1101 of 324 patients (30%) are still at risk. The analysis is based on the TCGA exome sample set. **B:**
1102 Kaplan-Meier survival curves for HNSC patients with and without *PRSS3* promoter mutations.
1103 The data is shown for a follow-up time of 2,000 days, at which point 42 out of 484 patients (9%)
1104 are still at risk. The analysis is based on the TCGA exome sample set. **C:** Genomic overview of
1105 *STK11*, zooming in on its combined splice sites region. The phyloP track shows evolutionary
1106 conservation of positions. **D:** A two-sided Wilcoxon rank sum test is performed for LUAD
1107 samples from the TCGA exome sample set, to investigate if samples mutated in the splice site
1108 region of *STK11* have a different gene expression level than samples without such mutations.
1109 **E:** Kaplan-Meier survival curves for LUAD samples with and without *STK11* splice site
1110 mutations. The data is shown for a follow-up time of 2,000 days, at which point 36 out of 438
1111 patients (8%) are still at risk. The analysis is based on the TCGA exome sample set.

1112 Supplementary Figure legends

1113 Figure 1-figure supplement 1

1114 The average number of mutations observed per sample per bp for each of the considered
1115 element types, as well as for intergenic regions.

1116 Figure 4-figure supplement 1

1117 Analysis of protein-coding genes to evaluate ncdDetect performance for scores defined by log-
1118 likelihoods and the number of mutations. The final null model is obtained by forward model
1119 selection. The QQ-plots show the p-values of all genes ($n = 19,256$) plotted against their
1120 uniform expectation under the null for each of the models considered. **A:** Results using log-
1121 likelihoods. **B:** Results using number of mutations. The corresponding plot for conservation
1122 scores are shown in Figure 4A.

1123 Figure 4-figure supplement 2

1124 The p-values (based on conservation scores) plotted as a function of the total number of
1125 mutations across samples observed per bp for all protein-coding genes. The point size indicates
1126 gene length. The mean number of mutations per bp is on average eight times higher for the
1127 COSMIC genes detected by ncdDetect compared to the undetected COSMIC genes.

1128 Figure 5-figure supplement 1

1129 Q-values and top-ten ranking elements for each of the three proposed scoring schemes.
1130 Protein-coding COSMIC genes, or non-coding elements associated to COSMIC genes, are
1131 highlighted in red. For each element, the region size is given together with the observed number
1132 of mutations and the expected number of mutations under the null model. **A:** The QQ-plot
1133 shows the p-values for all protein-coding genes (n=19,256) plotted against their uniform
1134 expectation under the null. 64 protein-coding genes are found to be significant (conservation
1135 scores). **B:** QQ-plot of p-values for all 5' UTRs (n=18,220). In total, 86 5' UTRs are significant.
1136 **C:** QQ-plot of p-values for all 3' UTRs (n=18,481), of which 16 are found to be significant. The
1137 complete sets of significant elements for each region type are given in Supplementary Files 1-3.
1138 Similar plots for promoter elements and splice sites are shown in Figure 5.
1139

1140 Figure 5-figure supplement 2

1141 The number of elements called significant for each of the three proposed scoring schemes, for
1142 each of the defined element types. The use of log-likelihoods results in the highest number of
1143 elements called significant across most element types, and the use of the number of mutations
1144 results in the fewest.

1145 Figure 6-figure supplement 1

1146 Examples of correlation between the uracil-DNA glycosylase deficiency signature statistic and
1147 *SMUG1* gene expression (first column), *UNG* gene expression (second column) and the product
1148 of *SMUG1* and *UNG* gene expression (third column) using TCGA exome data for seven
1149 different cancer types (rows). The correlation is assessed using one-sided Spearman's
1150 correlation tests. For some cancer types, the correlation coefficients are positive, though these
1151 cases are not significant and generally based on few samples.

1152 Figure 5-figure supplement 3

1153 Length distributions of all defined element types.

1154 Figure 3-figure supplement 1

1155 Illustration of time complexity of the ncdDetect algorithm. Each point illustrates the CPU time in
1156 seconds used to calculate the background score distribution for a candidate region of a given
1157 length for a given number of samples.

1158

1159 Tables

1160 Table 1

1161

element type	number of elements	median element length (bps)	percentage of genome covered
protein-coding genes	20,153	1,296	1.19
promoter elements	20,052	848	0.69

splice sites	18,682	30	0.03
3' UTRs	19,346	1,007	1.06
5' UTRs	19,078	259	0.25

1162

1163 Overview of elements analysed with ncdDetect. Regions located on chromosome X and Y are
1164 excluded from the analyses (Material and methods: Candidate elements).

1165 Source data legends

1166 Figure 1-source data 1

1167 The number of mutations observed for each of the 505 samples. This data set relates to Figure
1168 1B.

1169 Figure 4-source data 1

1170 P-values obtained on protein-coding genes for each of the five models considered. The p-values
1171 are obtained using conservation scores. This data set relates to Figure 4A.

1172 Figure 4-source data 2

1173 COSMIC Gene Census recall data. The fraction of recalled COSMIC genes in the top ncdDetect
1174 candidates, for the three different scoring schemes. This data set relates to Figure 4C.

1175 Figure 4-source data 3

1176 The background score distribution for the protein-coding gene TP53 obtained with conservation
1177 scores. This data set relates to Figure 4D.

1178 Figure 4-source data 4

1179 The background score distribution for the protein-coding gene PIK3CA obtained with
1180 conservation scores. This data set relates to Figure 4E.

1181 Figure 4-source data 5

1182 The background score distribution for the protein-coding gene SLFN11 obtained with
1183 conservation scores. This data set relates to Figure 4F.

1184 Figure 5-source data 1

1185 P-values obtained on promoters and splice sites using conservation scores. This data set
1186 relates to Figure 5A-B. P-values obtained using log-likelihoods and number of mutations as
1187 scores are in Supplementary Files 2-3.

1188 Figure 6-source data 1

1189 The defined uracil-DNA glycosylase deficiency signature statistic for each sample of the cancer
1190 types GBM, BLCA, CRC, BRCA, LUAD, SKCM and UCEC. For each sample, the SMUG1
1191 mutation status is indicated. This data set relates to Figure 6C.

1192 Figure 6-source data 2

1193 The defined uracil-DNA glycosylase deficiency signature statistic, as well as SMUG1 gene
1194 expression, UNG gene expression, and SMUG1xUNG gene expression for TCGA exome
1195 samples. Expression values are RSEM values. This data set relates to Figure 6D as well as
1196 Figure 6-figure supplement 1.

1197 Figure 7-source data 1

1198 STK11 mutation status and STK11 gene expression (RSEM) for 469 LUAD TCGA exome
1199 samples. This data set relates to Figure 7D.

1200 Figure 1-figure supplement 1-source data 1

1201 The number of mutations per sample per bp for the defined element types. This data set relates
1202 to Figure S1.

1203 Figure 4-figure supplement 1-source data 1

1204 P-values obtained on protein-coding genes for each of the models considered. The p-values are
1205 obtained using log-likelihoods and number of mutations as scores. This data set relates to
1206 Figure 4-figure supplement 1.

1207 **Figure 4-figure supplement 2-source data 1**

1208 For each protein-coding gene, the gene length, the number of observed mutations across all
1209 505 samples, and the p-value obtained using conservation scores are given. Furthermore, the
1210 COSMIC status is indicated. This data set relates to Figure 4-figure supplement 2.

1211 **Figure 5-figure supplement 1-source data 1**

1212 P-values obtained on protein-coding genes, 3' UTRs and 5' UTRs using conservation scores.
1213 This data set relates to Figure 5-figure supplement 1A-C. P-values obtained using log-
1214 likelihoods and number of mutations as scores are in Supplementary Files 2-3.

1215 **Figure 5-figure supplement 2-source data 1**

1216 The number of elements called significant for each of the three proposed scoring schemes, for
1217 each of the defined element types. This data set relates to Figure 5-figure supplement 2.

1218 **Figure 5-figure supplement 3-source data 1**

1219 The length of each of the analysed elements. This data set relates to Figure 5-figure
1220 supplement 3.

1221 **Figure A2-source data 1**

1222 P-value and gene length for each protein-coding gene. The p-values are obtained with and
1223 without the overdispersion-based rate adjustment. This data set relates to Figure A2.

1224 **Figure A3-source data 1**

1225 COSMIC Gene Census recall data. The fraction of recalled COSMIC genes in the top ncdDetect
1226 and ExInAator candidates. This data set relates to Figure A3.
1227

1228 **Supplementary files**

1229 **Supplementary File 1**

1230 File name: Supplementary File 1.xlsx.

1231 Title of data: Supplementary results.

1232 Description of data: P-values obtained with ncdDetect using conservation scores.

1233 Supplementary File 2

1234 File name: Supplementary File 2.xlsx.

1235 Title of data: Supplementary results.

1236 Description of data: P-values obtained with ncdDetect using log likelihood scores.

1237 Supplementary File 3

1238 File name: Supplementary File 3.xlsx.

1239 Title of data: Supplementary results.

1240 Description of data: P-values obtained with ncdDetect using number of mutations as scores.

1241 Supplementary File 4

1242 File name: Supplementary File 4.xlsx.

1243 Title of data: Region definitions.

1244 Description of data: Definitions of candidate elements: Promoter regions, splice sites, 3' UTRs,

1245 5' UTRs and protein-coding genes.

1246 Supplementary File 5

1247 File name: Supplementary File 5.xlsx.

1248 Title of data: Expression analyses and a uracil-DNA glycosylase deficiency signature statistic.

1249 Description of data: General expression analyses results and analyses performed to investigate

1250 the impact of SMUG1 mutations on expression levels as well as the uracil-DNA glycosylase

1251 deficiency signature statistic.

1252 **Supplementary File 6**

1253 File name: Supplementary File 6.xlsx.

1254 Title of data: Correlation between mutation status and survival.

1255 Description of data: Data overview and results obtained from the survival analyses.

1256 **Supplementary File 7**

1257 File name: Supplementary File 7.xlsx.

1258 Title of data: Information on how to access expression and survival TCGA data sets.

1259 Description of data: Overview of the specific TCGA samples included in the expression and
1260 survival analyses.

1261

1262

1263

1264

1 Appendix

2 Overdispersion-based rate adjustment

3 As our understanding of the mutational process is limited and as we do not know all relevant
4 explanatory variables for all our samples, there will always be a difference between the
5 predicted and actual mutation rate (Figure A1). The unaccounted for explanatory variables are
6 likely to have auto-correlated regional effects. The effect of differences between actual and
7 predicted mutation rates will thus be accumulated along elements and be most pronounced for
8 long elements (Figure A2). Even small biases in the predicted versus actual mutation rate may
9 become significant if elements are sufficiently long. In our case, the protein-coding genes are
10 the longest element type and therefore the most likely to be affected by such biases.

11
12 The unavoidable difference between actual and predicted mutation rates across elements and
13 samples will increase the unexplained variance and lead to an overdispersion of the number of
14 mutations per element (or other test statistics based on it). By capturing and taking this
15 overdispersion into account, the specificity of the method can be improved, though not the
16 power, which depends on reducing the unexplained variance by better mutational null models.
17 To correct for overdispersion, we adjust each sample- and position-specific mutational
18 probability by an overdispersion-based mutation rate correction factor.

19
20 The overdispersion-based rate adjustment is modelled with a beta binomial model. For a region
21 of length L_i having X_i mutations, we have

$$22 \quad X_i \sim \text{Binomial}(L_i \cdot N_i, p_i),$$
$$p_i \sim \text{Beta}(\alpha_i, \beta_i).$$

23 The parameters α_i and β_i are constrained to satisfy that the expected mutation probability
 24 $\mu_i = \mathbb{E}[p_i] = \frac{\alpha_i}{\alpha_i + \beta_i}$ equals \hat{p}_i , the average mutation rate in the region predicted by the logistic
 25 regression model. Left with one degree of freedom the overdispersion is modelled with the
 26 parameter $\gamma = \frac{\text{SD}(p)}{\mathbb{E}(p)} = \sqrt{\frac{\beta}{\alpha} \frac{1}{\alpha + \beta + 1}}$. We can express α and β in terms of μ and γ :

$$\alpha = \frac{1 - \mu - \gamma^2 \mu}{\gamma^2},$$

$$\beta = \frac{1 - \mu}{\mu} \cdot \frac{1 - \mu - \gamma^2 \mu}{\gamma^2}.$$

27 In this alternative parameterization our model becomes

$$X_i \sim \text{Binomial}(L_i \cdot N_i, p_i),$$

$$p_i \sim \text{Beta}(\hat{p}_i, \gamma).$$

28 The parameter γ is shared across all regions, and we estimate it by numerically maximizing the
 29 likelihood function

$$L(\gamma) = \prod_{i=1}^K \binom{NL_i}{X_i} \frac{\text{B}(X_i + \alpha, NL_i - X_i + \beta)}{\text{B}(\alpha, \beta)}$$

$$= \prod_{i=1}^K \binom{NL_i}{X_i} \frac{\text{B}(X_i + \frac{1 - \mu - \gamma^2 \mu}{\gamma^2}, NL_i - X_i + \frac{1 - \mu}{\mu} \cdot \frac{1 - \mu - \gamma^2 \mu}{\gamma^2})}{\text{B}(\frac{1 - \mu - \gamma^2 \mu}{\gamma^2}, \frac{1 - \mu}{\mu} \cdot \frac{1 - \mu - \gamma^2 \mu}{\gamma^2})},$$

30 where B is the beta function. To avoid that regions under positive selection affect the estimate of
 31 γ , we filter out the top 5% and bottom 5% of regions, where the observed number of mutations
 32 deviates the most from the expected number of mutations. For protein-coding genes we further
 33 explicitly filter out COSMIC genes.

34

35 Inspecting the QQ-plot of the p-values for protein-coding genes shows that the overdispersion-
 36 based rate adjustment improves the overall fit of the p-values to uniformity under the null and

37 reduces the inflation of the tail of the distribution (Figure A2A). For short genes (<700 bp), the
38 QQ-plots show a near perfect fit of the p-values to the uniform expectation with known or likely
39 cancer drivers standing out (Figure A2B). For long genes (>3,000 bp), the fit is much improved
40 by the rate-adjustment, though some inflation is still present, resulting in significant calls that are
41 likely false positives (e.g., *MUC4*, *PLIN4*, etc.) (Figure A2C).

42 ncdDetect compared to other non-coding cancer driver detection 43 methods

44 In order to benchmark the performance of ncdDetect, we compared our results to those
45 obtained with two other non-coding cancer driver detection methods, ExInAator (Lanzós et al.
46 2017) and LARVA (Lochovsky et al. 2015). ExInAator is designed for the analysis of lncRNAs,
47 but is also applicable on protein-coding genes. We thus compared ncdDetect and ExInAator on
48 protein-coding genes. At the time of writing, LARVA does not support discontinuous element
49 types (e.g. the joint analysis of multiple exons within a single gene). The promoter elements
50 analysed in the present paper are in principle contiguous. However, as we subtract overlapping
51 annotations from 5' UTRs, they can be discontinuous in some cases. We have run LARVA on
52 our promoter definitions, without removing any overlapping annotations from other element
53 types.

54

55 The benchmarking on protein-coding genes is performed using the COSMIC Cancer Gene
56 Census as a true positive set (Forbes et al. 2014). The promoter benchmarking is not as
57 straightforward, given the lack of a true positive set for this element type. We compare the
58 results using previously described promoter cancer driver candidates.

59 Performance on protein-coding genes

60 Where ncdDetect has a tendency to suffer from a significant fraction of false positive
61 predictions, ExInAator has a tendency to suffer from a significant fraction of false negative
62 predictions. The published ExInAator results on protein-coding genes, based on the same 505
63 cancer samples analysed here, contain p-values for 19,309 protein-coding genes, where three
64 are significant ($q < 0.10$). One of those three significant genes is in the COSMIC database. For
65 ncdDetect, 64 genes are called significant ($q < 0.10$), of which 15 ($\approx 23\%$) are COSMIC genes.
66 Notably, six of the top-ten ncdDetect candidates are COSMIC genes ($p = 2.0 \cdot 10^{-7}$, Fisher's
67 exact test). This is the case for three of the top-ten ExInAator candidates ($p = 3.3 \cdot 10^{-3}$, Fisher's
68 exact test). In general, ExInAator predicts much fewer candidates than ncdDetect, and thus have
69 a lower false positive rate. However, ncdDetect has a higher COSMIC Census recall rate, i.e. it
70 performs better at ranking genes compared to ExInAator (Figure A3).

71
72 Looking closer at the top-15 ncdDetect protein-coding candidates, we find that nine are
73 COSMIC genes (Table A1). Non-COSMIC genes in the list include *PRSS3*, *KRTAP9-1*,
74 *KRTAP4-5* and *BCLAF1*. All of these genes have some reported cancer association: The
75 expression of *PRSS3* has been shown to be upregulated in metastatic prostate cancer and is
76 also associated to pancreatic and lung cancer (Jiang et al. 2010; Hockla et al. 2012; Marsit et al.
77 2005). Although The Keratin associated proteins *KRTAP9-1* and *KRTAP4-5* have no wide
78 spread reported role in cancer, a recent study found that they can play a part in malignant
79 progression (Berens et al. 2017). Finally, *BCLAF* has been associated to colon cancer (Zhou et
80 al. 2014). The remaining two genes in the list, *MUC4* and *AL390778.1* have no reported cancer
81 driver potential. Interestingly *MUC4* continues to be significant in our analyses, even after
82 overdispersion-based rate adjustment.

83 Performance on non-coding regulatory elements

84 The LARVA analysis of the 20,052 defined promoter elements yields 16 significant candidates
85 ($q < 0.10$). ncdDetect agrees, and calls all of these 16 promoters significant, along with an
86 additional 144 candidates ($q < 0.10$). Several of the cases detected by ncdDetect, and not by
87 LARVA, have previously been described to be associated with cancer. These include *PLEKHS1*
88 and *WDR74* as described in the main text. Further, promoter mutations in *DPH3* and *OXNAD1*
89 have been associated to skin cancers (Denisova et al. 2015). A number of the ncdDetect
90 identified candidates are also identified in an earlier cancer study (Weinhold et al. 2014),
91 including *SMUG1* (Figure A4). Finally, the protein-coding genes associated to the promoters
92 *KDM5A*, *CNOT3* and *NCOR1* are COSMIC genes, and detected solely by ncdDetect. The
93 promoter region of *PRSS3*, a case study in the main text, is also detected by ncdDetect alone.
94 Taken together, several of the ncdDetect promoter candidates that are not detected by LARVA
95 have previously reported cancer driver potential.

97 Algorithmic details of ncdDetect

98 Assume that a candidate element has m positions and that somatic mutations are called for a
99 total of k samples. The sample- and position-specific probabilities of a mutation are predicted
100 using the null model. The four outcomes of the model are one type of transition (TS =
101 $TS_{\{A \rightarrow G, G \rightarrow A\}}$), two types of transversions ($TV_1 = TV_{\{A \rightarrow T, G \rightarrow T\}}$ and $TV_2 = TV_{\{A \rightarrow C, G \rightarrow C\}}$), as well as
102 the reference class of no mutation (NM). The corresponding probabilities for the i 'th sample
103 ($i = 1, \dots, k$) and position j ($j = 1, \dots, m$) are (Figure 3C)

$$104 \left(\pi_{ij}^{\text{TS}}, \pi_{ij}^{\text{TV}_1}, \pi_{ij}^{\text{TV}_2}, \pi_{ij}^{\text{NM}} \right). \quad (1)$$

105 Associated to each outcome is a sample- and position-specific score (Figure 3B)

$$\left(\text{score}_{ij}^{\text{TS}}, \text{score}_{ij}^{\text{TV}_1}, \text{score}_{ij}^{\text{TV}_2}, \text{score}_{ij}^{\text{NM}}\right).$$

106 Let $\text{obs}(i, j)$ indicate the observed outcome for position i and sample j (Figure 3D). Then the
 107 observed score for position i and sample j is $\text{score}_{ij}^{\text{obs}(i,j)}$, the cumulated observed sample-
 108 specific score is given by (Figure 3E)

$$\text{score}_i^{\text{obs}} = \sum_{j=1}^m \text{score}_{ij}^{\text{obs}(i,j)},$$

109 and the overall score is (Figure 3F)

$$\text{score}_{\text{overall}}^{\text{obs}} = \sum_{i=1}^k \text{score}_i^{\text{obs}} = \sum_{i=1}^k \sum_{j=1}^m \text{score}_{ij}^{\text{obs}(i,j)}.$$

110 We now describe how to determine the null distribution for the overall score. First consider the
 111 null distribution for the sample-specific score (Figure 3E). Let $Z(i, j)$ be the stochastic variable
 112 that indicates the outcome for position i and sample j . Each of the four outcomes (TS , TV_1 , TV_2 ,
 113 NM) happen with the probability determined by equation (1). The cumulated sample-specific
 114 score distribution is thus the distribution of the stochastic variable

$$A_{i,m} = \sum_{j=1}^m \sum_{z \in \{\text{TS}, \text{TV}_1, \text{TV}_2, \text{NM}\}} \text{score}_{ij}^z \mathbf{1}(Z(i, j) = z),$$

115 where $\mathbf{1}(\cdot)$ is the indicator function. If we assume that scores are non-negative integers we have
 116 the recursion from one position to the next

$$P(A_{i,j} = s) = \sum_{\ell=0}^s P(A_{i,(j-1)} = \ell) \sum_{z \in \{\text{TS}, \text{TV}_1, \text{TV}_2, \text{NM}\}} P(\text{score}_{i,j}^z = s - \ell).$$

117 Second consider the null distribution for the overall score (Figure 3G)

$$B_k = \sum_{i=1}^k A_{i,m}.$$

118 A similar recursion as before holds for the overall score distribution. We can include the next
 119 sample from the recursion

$$P(B_k = s) = \sum_{\ell=0}^s P(B_{k-1} = \ell)P(A_{k,m} = s - \ell).$$

120 The final p-value for the element of interest is

$$P(B_k \geq \text{score}_{\text{overall}}^{\text{obs}}) = 1 - P(B_k < \text{score}_{\text{overall}}^{\text{obs}}).$$

121 Figure legends

122 Figure A1

123

124 Illustration of the motivation behind the overdispersion-based rate adjustment. For candidate
125 element A, we overestimate the mutation rate, and thus end up with a conservative p-value for
126 this element when analysing it with ncdDetect. For candidate element B, on the other hand, we
127 underestimate the mutation rate. In this case ncdDetect will produce a p-value that is too small,
128 creating a potential false positive call. The effect of underestimating the mutation rate will be
129 greater for longer candidate elements.

130 Figure A2

131 QQ-plots of p-values obtained with and without the overdispersion-based rate adjustment. **A:**
132 QQ-plots of all protein-coding genes (excluding *TP53* for illustration purposes). **B:** QQ-plots of
133 protein-coding genes shorter than 700 bp. For the shorter genes, the p-values are not
134 particularly inflated. The overdispersion-based rate adjustment does not affect the distribution of
135 p-values much. **C:** QQ-plots of protein-coding genes longer than 3,000 kb. For the longer
136 genes, the p-values are inflated, and the overdispersion-based rate adjustment effectively
137 corrects for much of this inflation.

138 Figure A3

139 COSMIC Gene Census recall plot. The fraction of COSMIC genes recalled in the top ncdDetect
140 and ExInAator candidates.

141

142 Figure A4

143 Illustration of overlap between significant elements found by ncdDetect and other non-coding
144 cancer driver screens. Highlighted elements are mentioned in the text. **A:** Overlap of promoter
145 elements found to be significant with ncdDetect and LARVA, as well as promoter elements
146 previously described in a non-coding cancer driver screen (Weinhold et al. 2014). We note that
147 *TERT* and *PLEKHS1* are also detected by a second non-coding driver screen (Melton et al.
148 2015). **B:** Overlap between 3' UTRs detected by ncdDetect and 3' UTRs detected by a previous
149 study (Weinhold et al. 2014). **C:** Overlap between 5' UTRs detected by ncdDetect and 5' UTRs
150 detected by a previous study (Weinhold et al. 2014). We note, that out of the 863 whole
151 genomes analysed in (Weinhold et al. 2014), 356 are sequenced by the TCGA. These samples
152 appear to be a subset of the 505 TCGA samples analysed here. The data sets are thus not
153 completely independent.

154 Tables

155 Table A1

156

157

rank	gene name	q-value	size (bp)	COSMIC	conclusion
1	TP53	1.15 x 10 ⁻²³⁵	1,378	true	putative cancer driver gene
2	PIK3CA	2.82 x 10 ⁻⁴⁴	3,207	true	putative cancer driver gene
3	KRAS	4.93 x 10 ⁻³³	708	true	putative cancer driver gene
4	PTEN	3.70 x 10 ⁻²⁵	1,212	true	putative cancer driver gene
5	PRSS3	6.80 x 10 ⁻²⁰	1,056	false	reported cancer association
6	MUC4	1.26 x 10 ⁻¹⁶	16,239	false	likely false positive due to length
7	KRTAP9-1	1.77 x 10 ⁻¹⁴	770	false	reported cancer association
8	BRAF	3.25 x 10 ⁻¹²	2,301	true	putative cancer driver gene
9	IDH1	1.76 x 10 ⁻¹⁰	1,245	true	putative cancer driver gene
10	KRTAP4-5	6.11 x 10 ⁻¹⁰	546	false	reported cancer association
11	NRAS	2.20 x 10 ⁻⁸	570	true	putative cancer driver gene
12	AL390778.1	5.95 x 10 ⁻⁸	735	false	no reported cancer driver properties
13	NFE2L2	3.15 x 10 ⁻⁷	1,890	true	putative cancer driver gene
14	FBXW7	7.35 x 10 ⁻⁷	2,618	true	putative cancer driver gene
15	BCLAF1	7.92 x 10 ⁻⁶	2,763	false	reported cancer association

158

159 Analysis of the top 15 ncdDetect protein-coding candidates.

160

161

162 References

163 Berens, E. B., G. M. Sharif, M. O. Schmidt, G. Yan, C. W. Shuptrine, L. M. Weiner, E. Glasgow,
164 A. T. Riegel, and A. Wellstein. 2017. "Keratin-Associated Protein 5-5 Controls Cytoskeletal
165 Function and Cancer Cell Vascular Invasion." *Oncogene* 36 (5): 593–605.

166 Denisova, Evgeniya, Barbara Heidenreich, Eduardo Nagore, P. Sivaramakrishna Rachakonda,
167 Ismail Hosen, Ivana Akrap, Víctor Traves, et al. 2015. "Frequent DPH3 Promoter Mutations
168 in Skin Cancers." *Oncotarget* 6 (34): 35922–30.

169 Forbes, S. A., D. Beare, P. Gunasekaran, K. Leung, N. Bindal, H. Boutselakis, M. Ding, et al.
170 2014. "COSMIC: Exploring the World's Knowledge of Somatic Mutations in Human
171 Cancer." *Nucleic Acids Research* 43 (D1): D805–11.

172 Hockla, A., E. Miller, M. A. Salameh, J. A. Copland, D. C. Radisky, and E. S. Radisky. 2012.
173 "PRSS3/Mesotrypsin Is a Therapeutic Target for Metastatic Prostate Cancer." *Molecular*
174 *Cancer Research: MCR* 10 (12): 1555–66.

175 Jiang, Guozhong, Fengyu Cao, Guoping Ren, Dongling Gao, Vipul Bhakta, Yunhan Zhang, Hua
176 Cao, et al. 2010. "PRSS3 Promotes Tumour Growth and Metastasis of Human Pancreatic
177 Cancer." *Gut* 59 (11): 1535–44.

178 Lanzós, Andrés, Joana Carlevaro-Fita, Loris Mularoni, Ferran Reverter, Emilio Palumbo,
179 Roderic Guigó, and Rory Johnson. 2017. "Discovery of Cancer Driver Long Noncoding
180 RNAs across 1112 Tumour Genomes: New Candidates and Distinguishing Features."
181 *Scientific Reports* 7 (January): 41544.

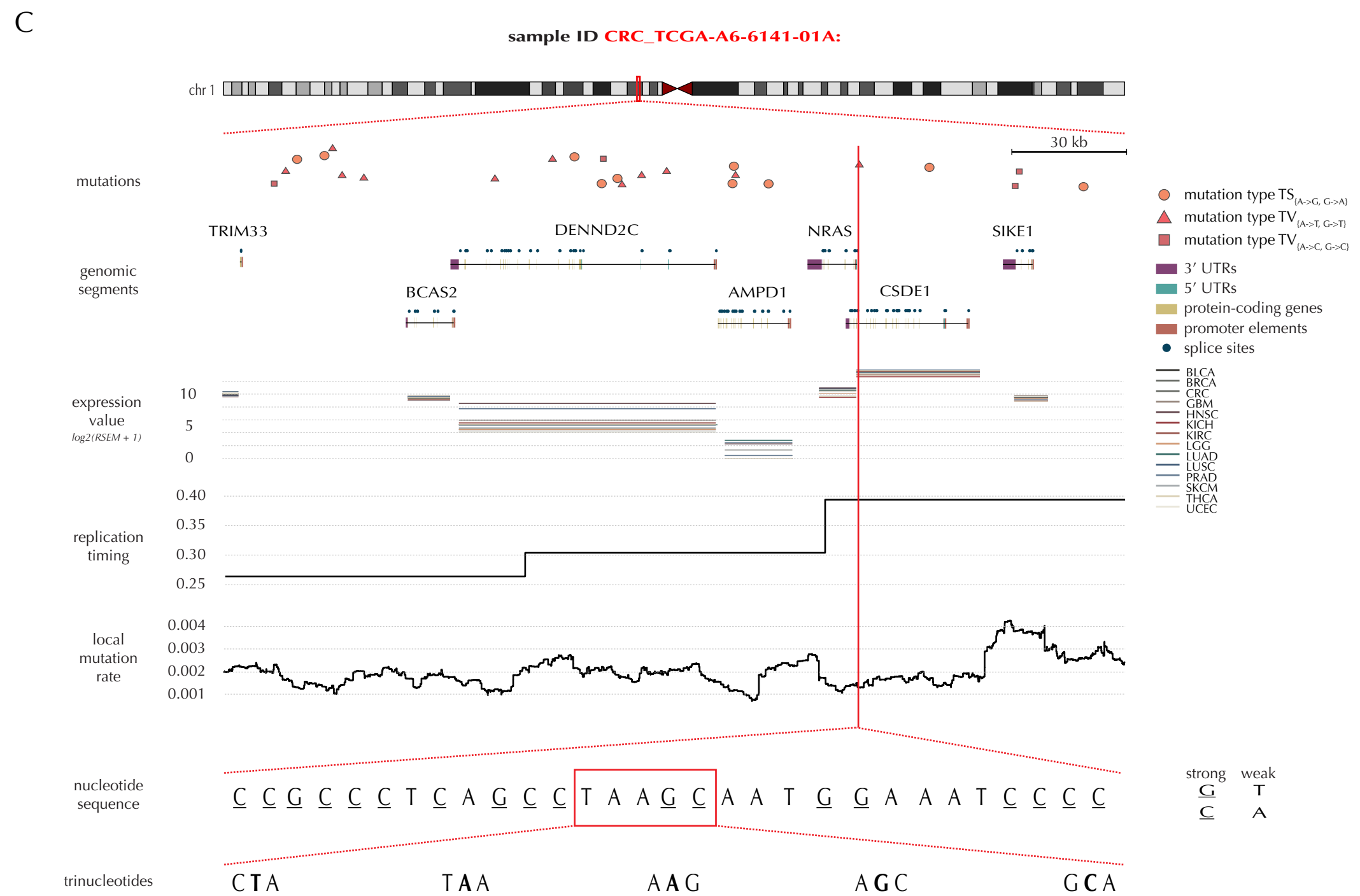
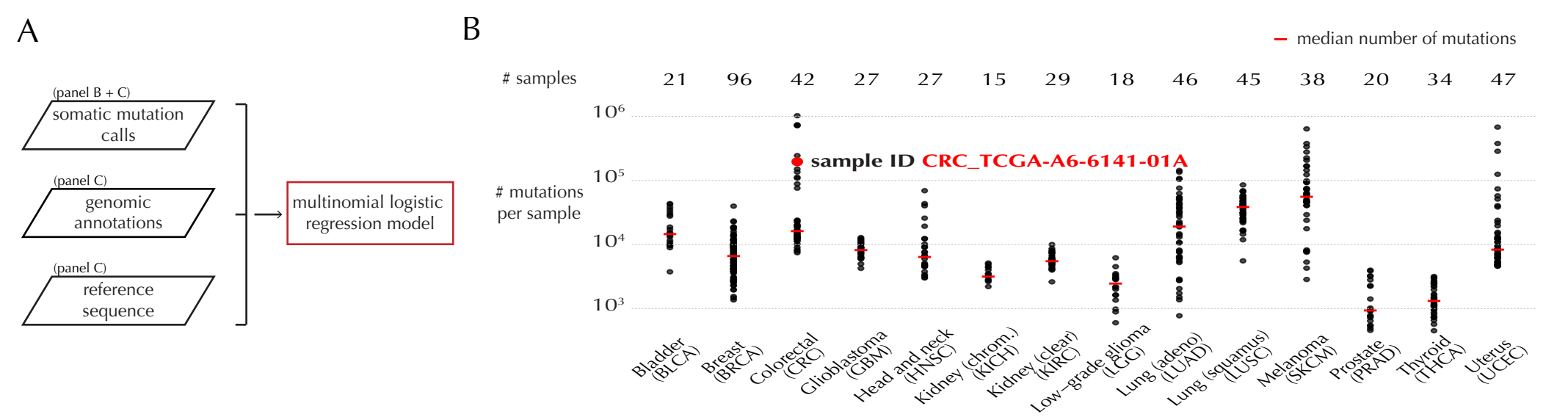
182 Lochovsky, Lucas, Jing Zhang, Yao Fu, Ekta Khurana, and Mark Gerstein. 2015. "LARVA: An
183 Integrative Framework for Large-Scale Analysis of Recurrent Variants in Noncoding
184 Annotations." *Nucleic Acids Research* 43 (17): 8123–34.

185 Marsit, Carmen J., Chinedu Okpukpara, Hadi Danaee, and Karl T. Kelsey. 2005. "Epigenetic
186 Silencing of the PRSS3 Putative Tumor Suppressor Gene in Non-Small Cell Lung Cancer."
187 *Molecular Carcinogenesis* 44 (2): 146–50.

188 Melton, Collin, Jason A. Reuter, Damek V. Spacek, and Michael Snyder. 2015. "Recurrent
189 Somatic Mutations in Regulatory Regions of Human Cancer Genomes." *Nature Genetics*
190 47 (7): 710–16.

191 Weinhold, Nils, Anders Jacobsen, Nikolaus Schultz, Chris Sander, and William Lee. 2014.
192 "Genome-Wide Analysis of Noncoding Regulatory Mutations in Cancer." *Nature Genetics*
193 46 (11): 1160–65.

194 Zhou, Xuexia, Xuebing Li, Yuanming Cheng, Wenwu Wu, Zhiqin Xie, Qiulei Xi, Jun Han,
195 Guohao Wu, Jing Fang, and Ying Feng. 2014. "BCLAF1 and Its Splicing Regulator SRSF10
196 Regulate the Tumorigenic Potential of Colon Cancer Cells." *Nature Communications* 5
197 (August): 4581.



A

(Figure 1 and panel B)

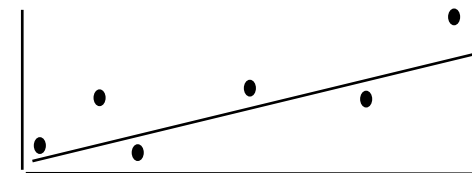
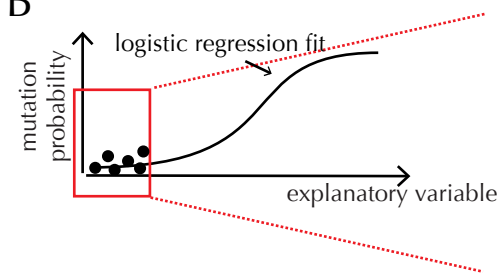
multinomial logistic regression model



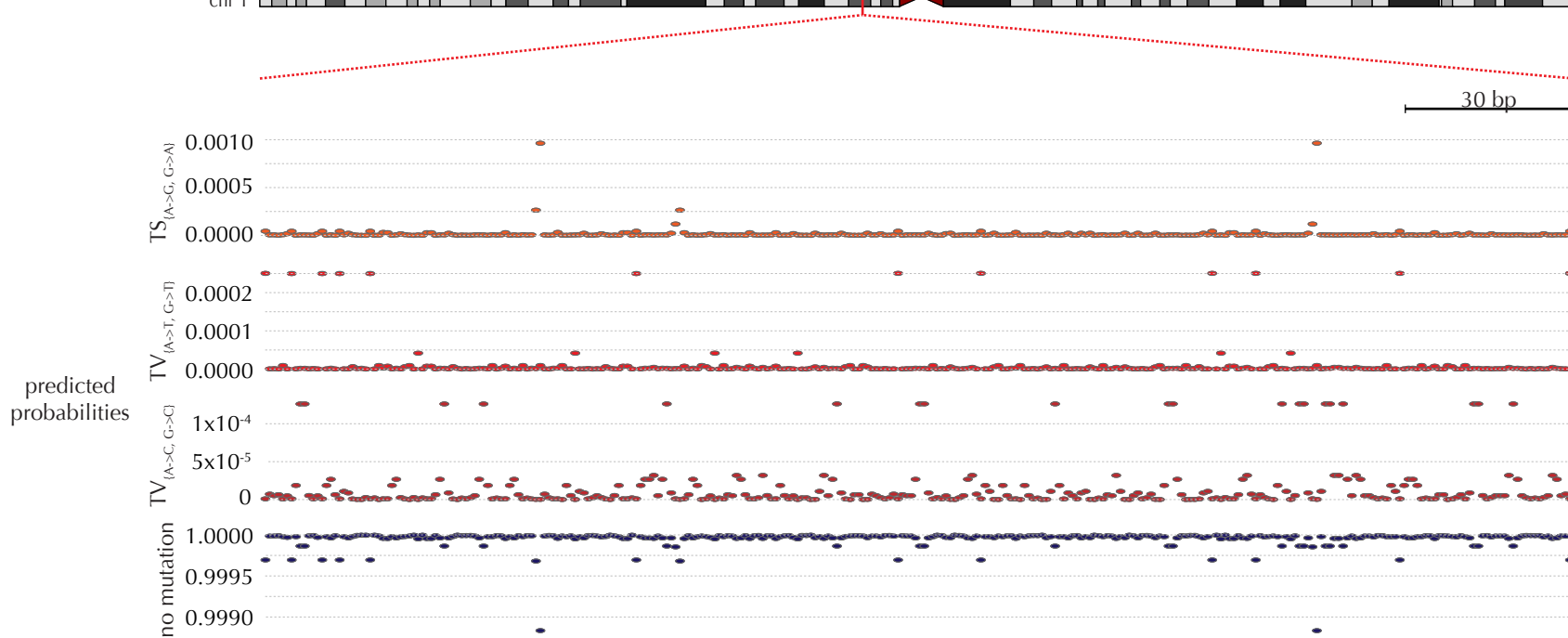
(panel C)

probabilities of mutation

B



C

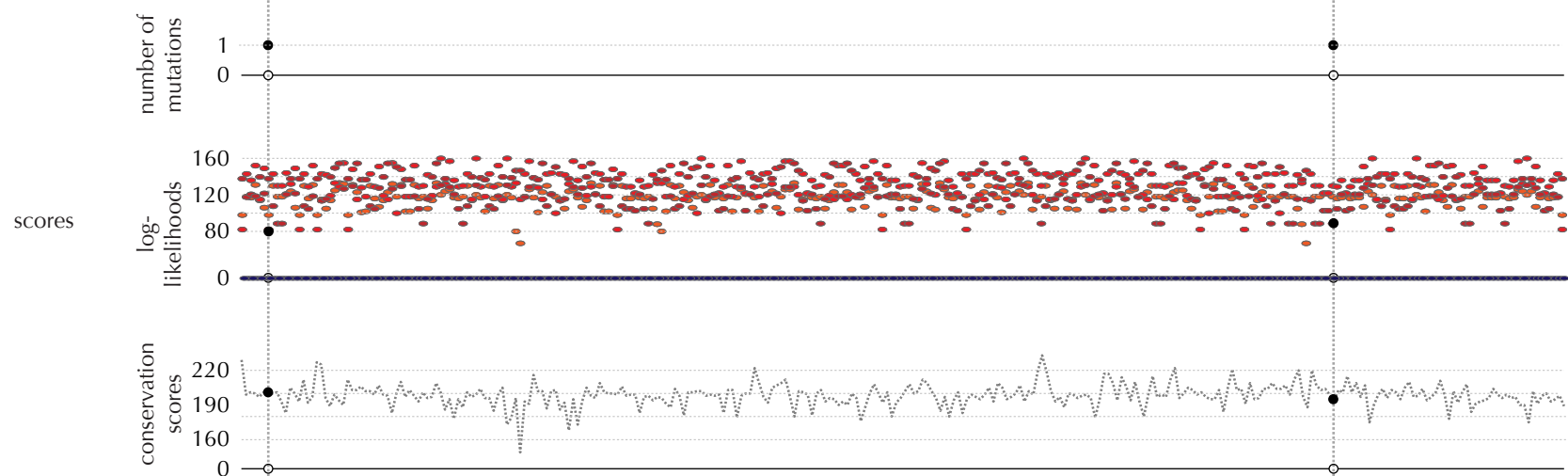
sample ID **CRC_TCGA-A6-6141-01A:**

D

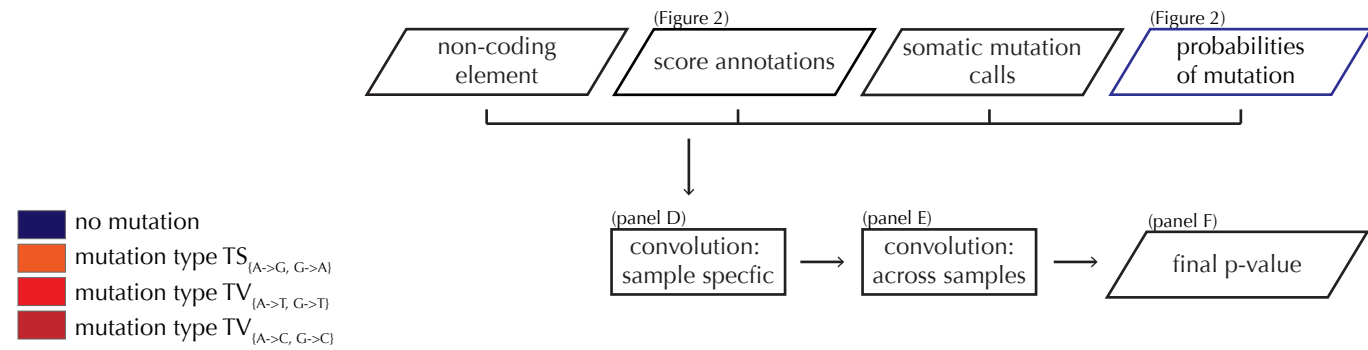
observed mutations



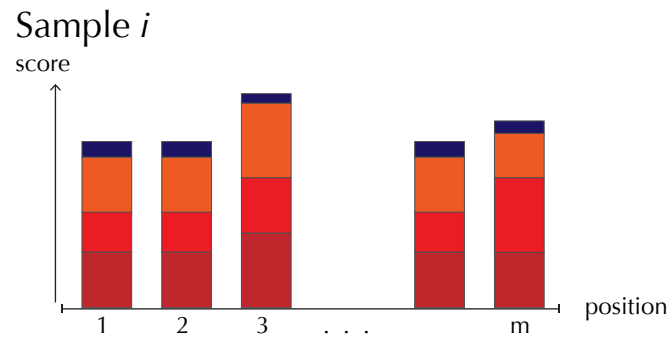
E



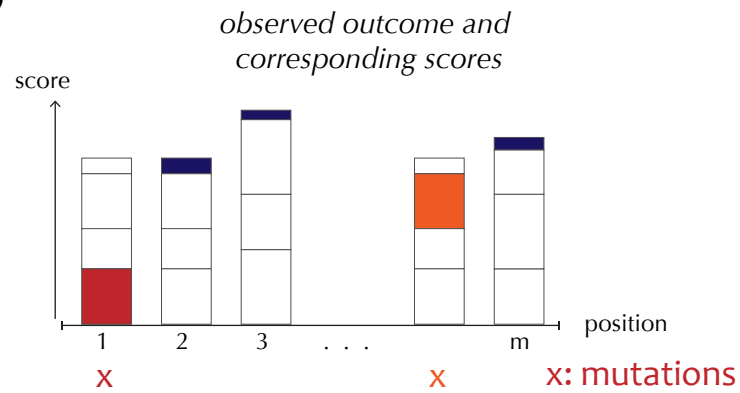
A



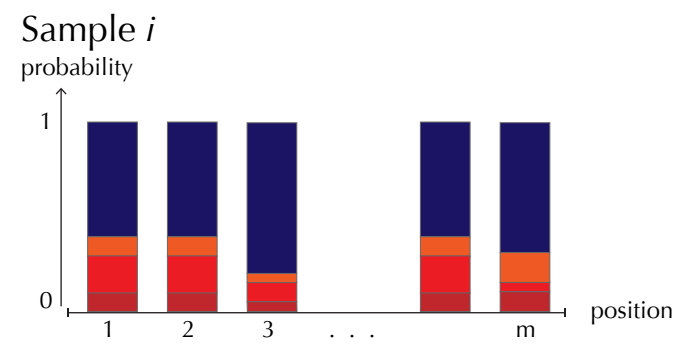
B



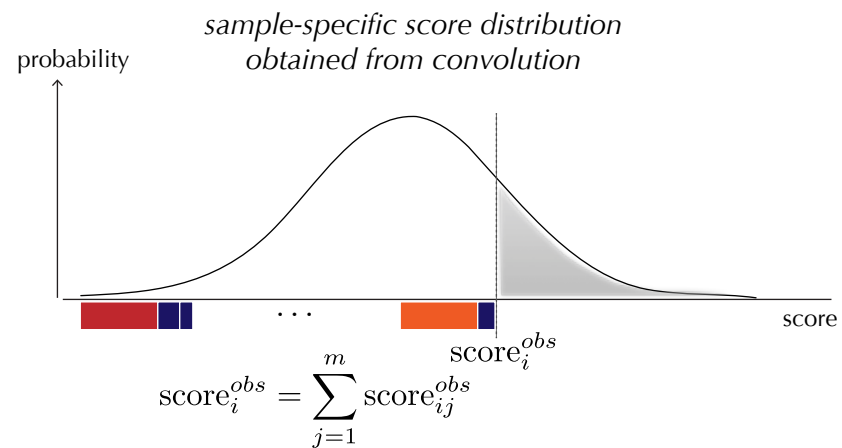
D



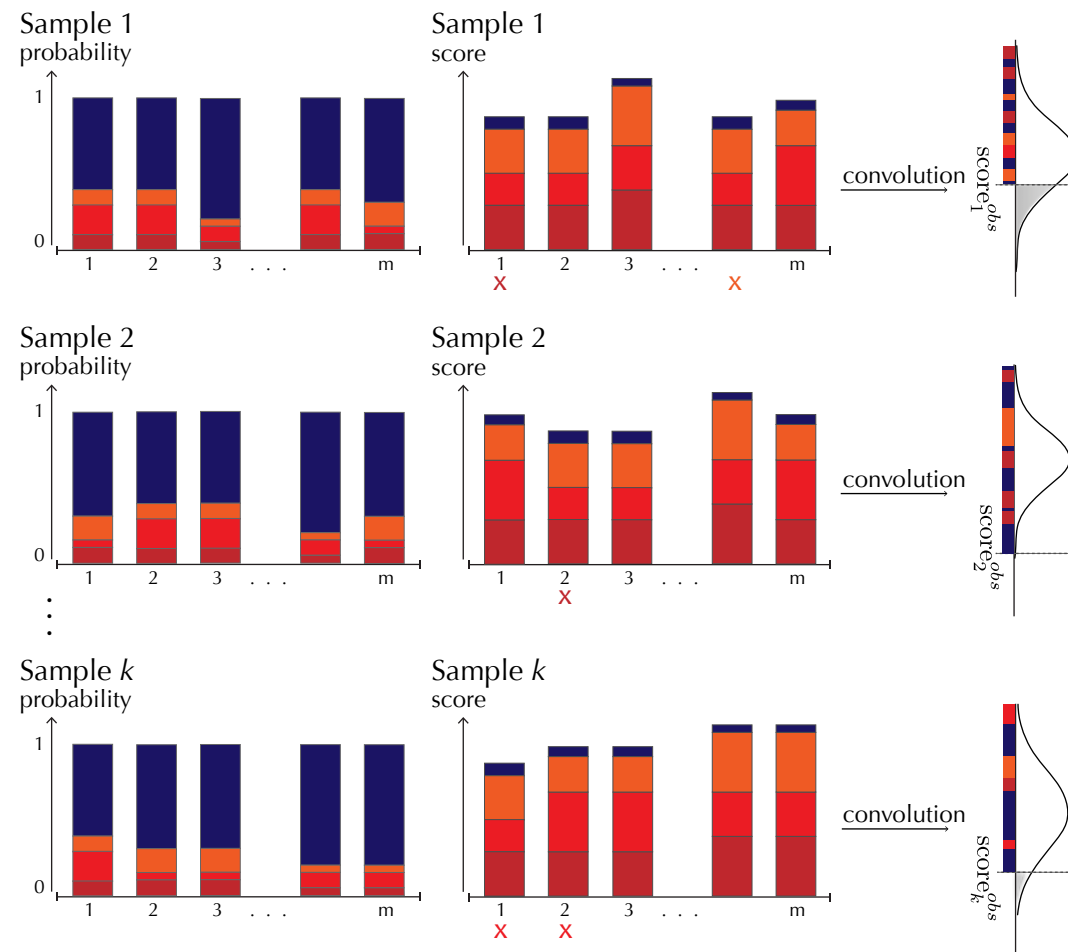
C



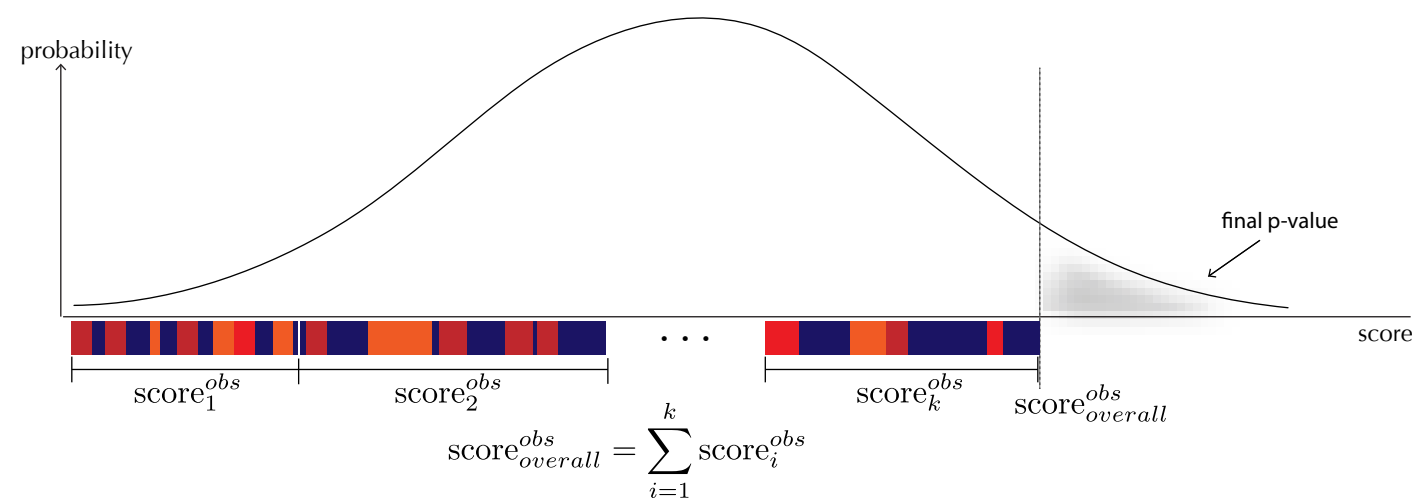
E

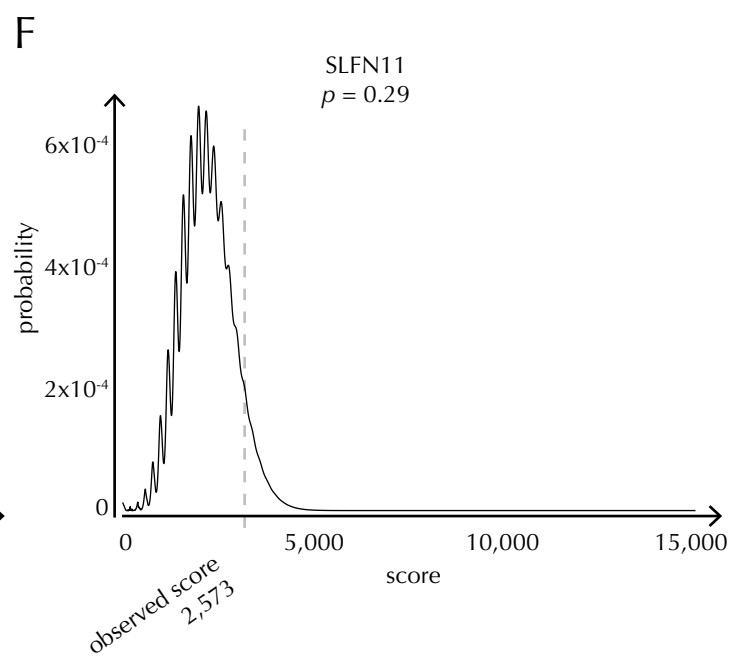
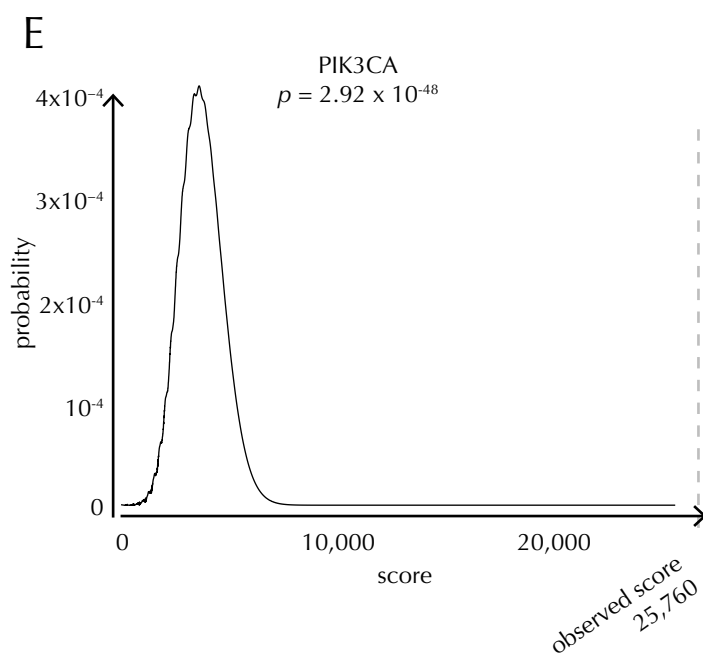
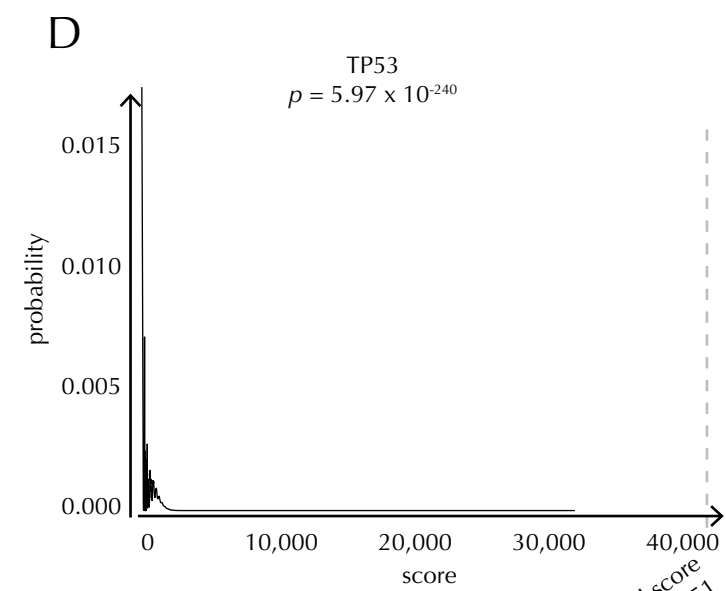
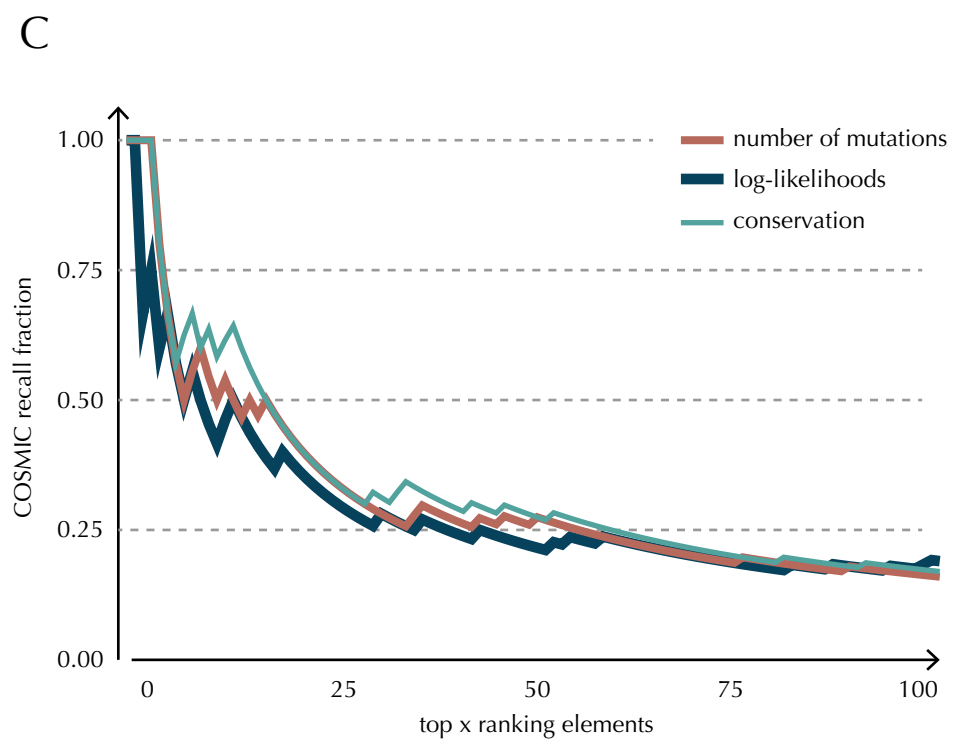
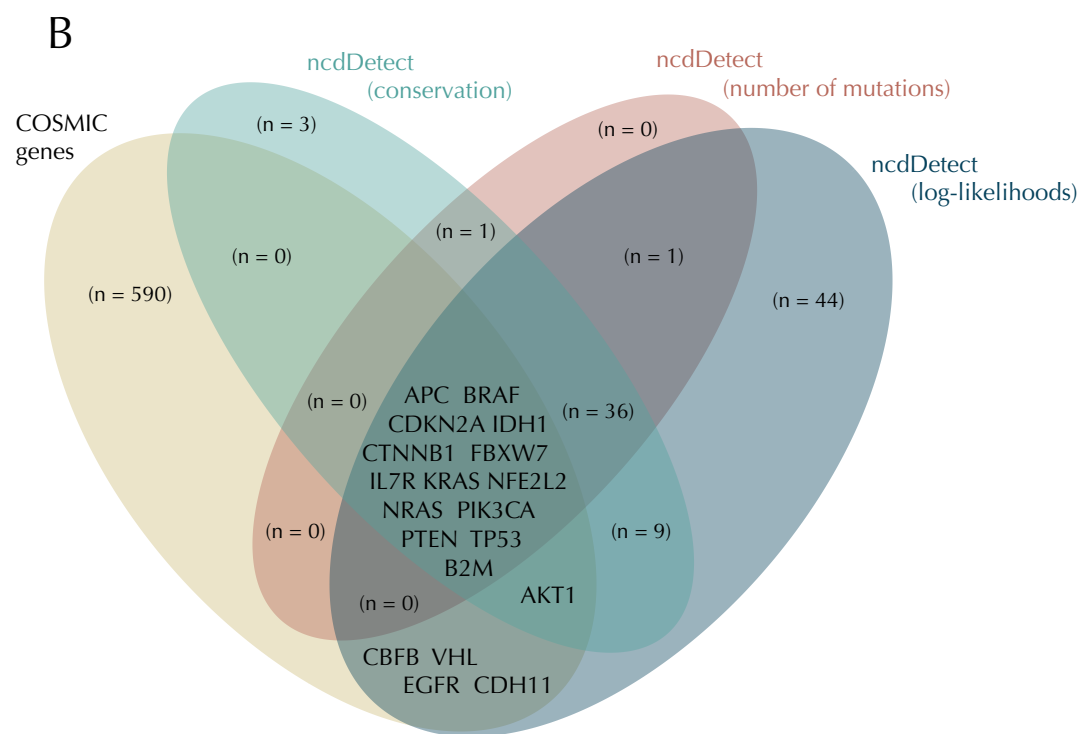
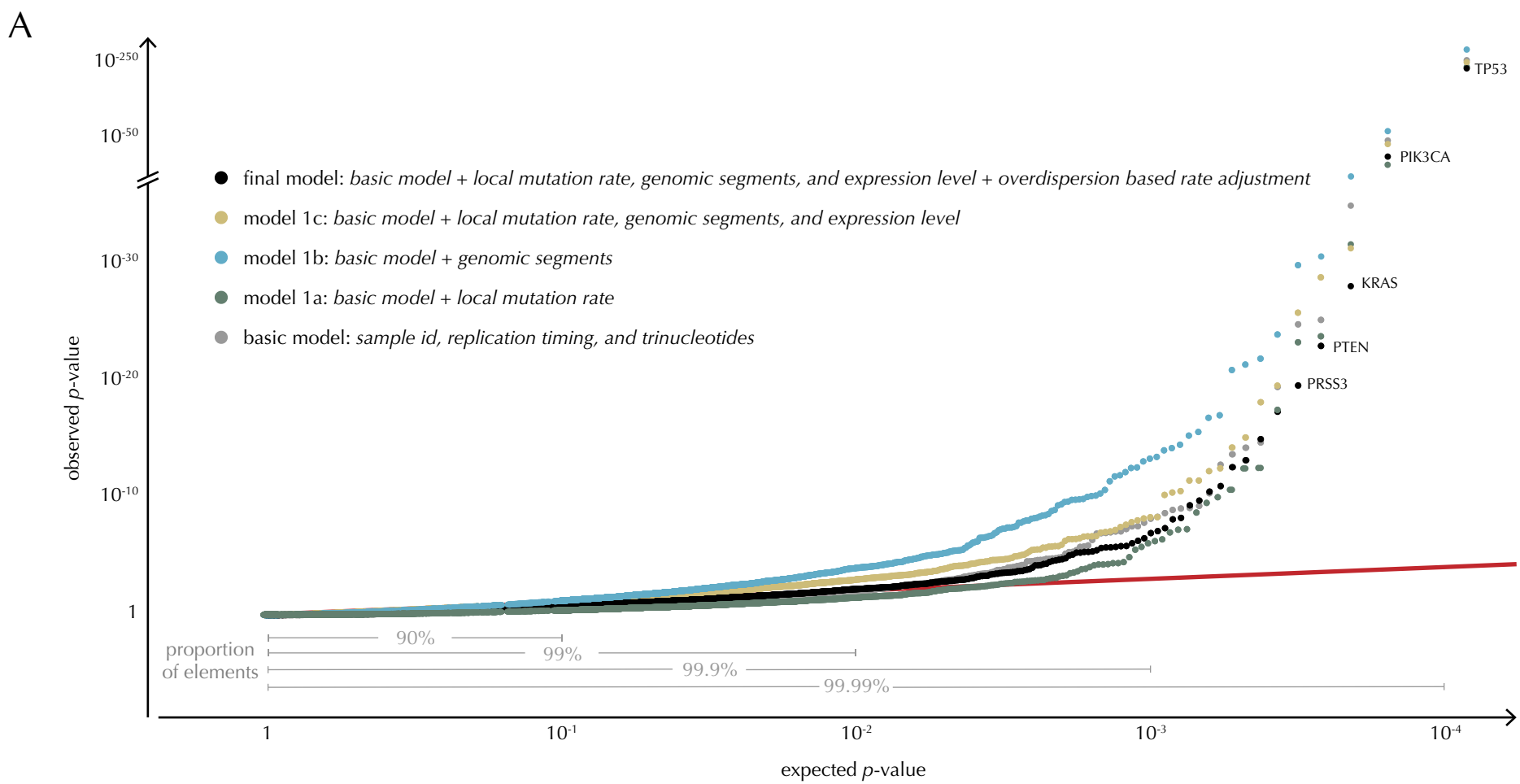


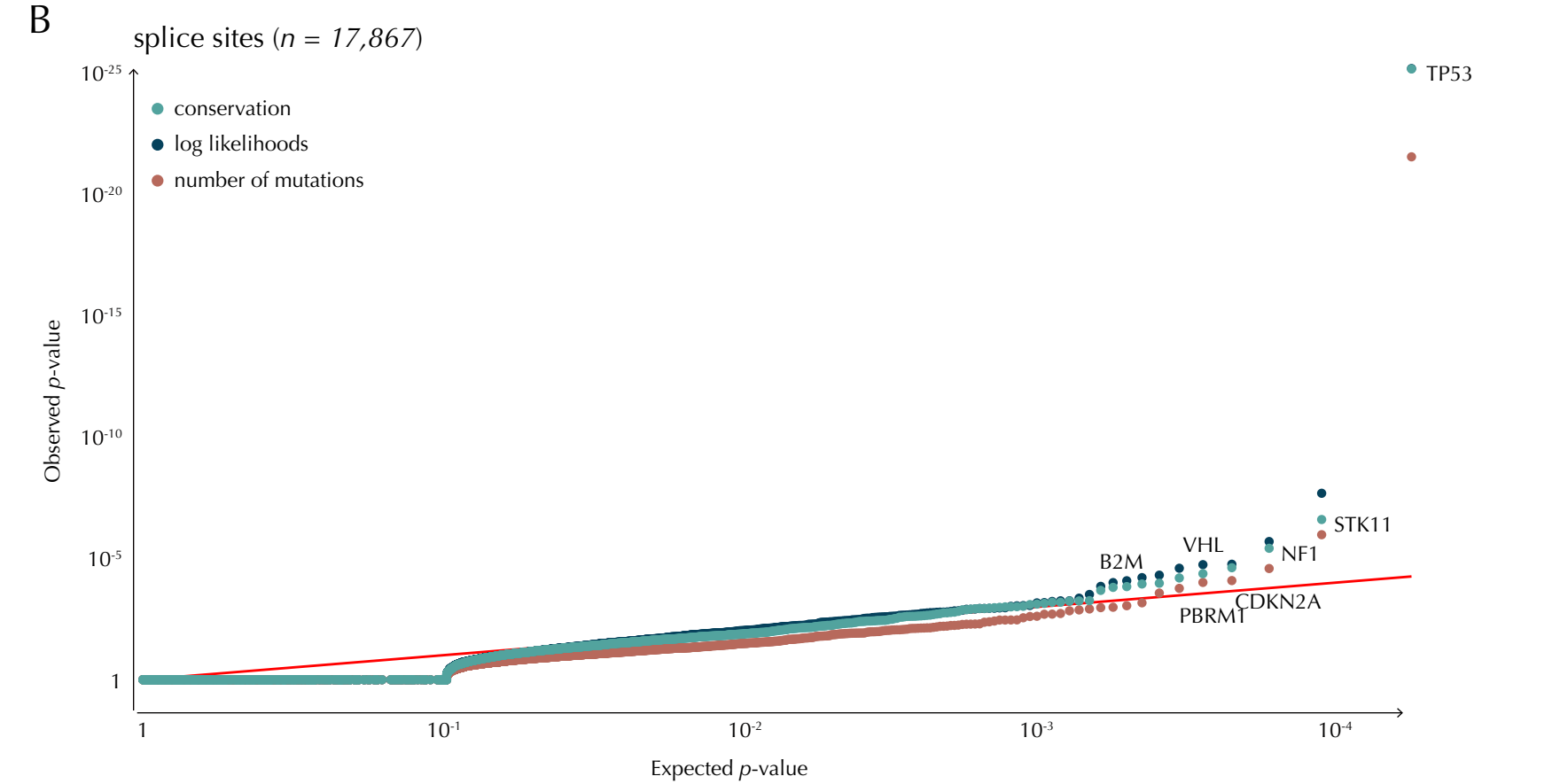
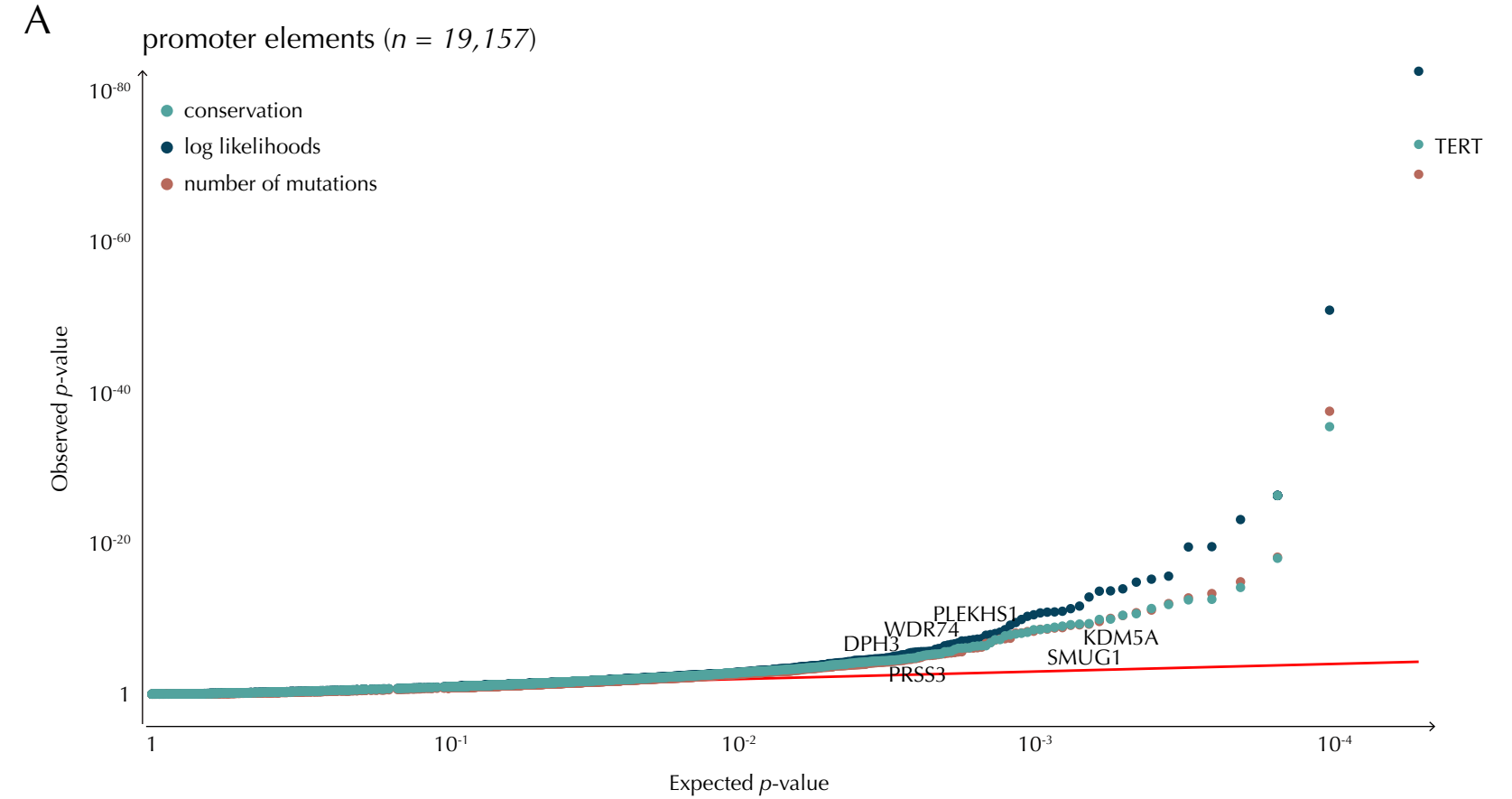
F



G







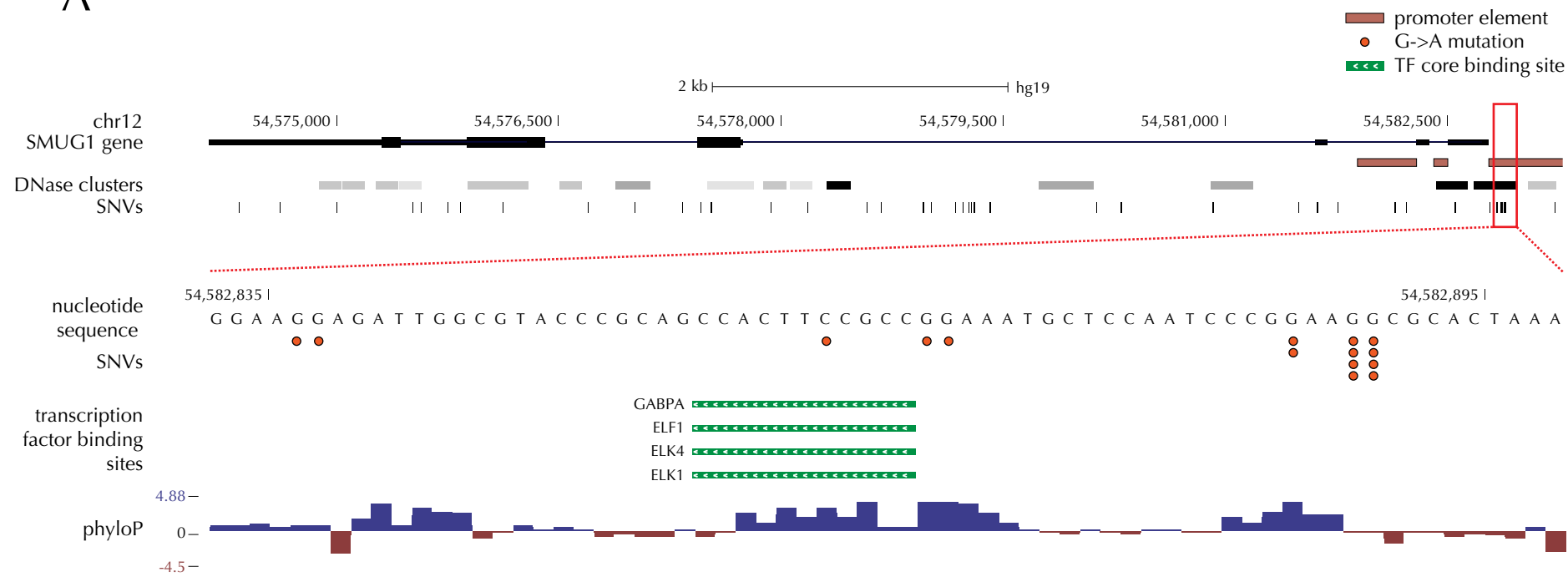
rank	conservation		log likelihoods		number of mutations	
	gene name	q-value	gene name	q-value	gene name	q-value
1	TERT	2.39×10^{-69}	TERT	4.79×10^{-79}	TERT	2.23×10^{-65}
2	FRG2B	3.28×10^{-32}	FRG2B	1.17×10^{-47}	FRG2B	2.94×10^{-34}
3	HLA-DRB5	6.19×10^{-15}	HLA-DRB5	2.93×10^{-23}	HLA-DRB5	4.40×10^{-15}
4	MUC3A	3.44×10^{-11}	MUC3A	3.42×10^{-20}	MUC3A	6.05×10^{-12}
5	AL645608.2	1.02×10^{-9}	SPATC1L	9.80×10^{-17}	AL645608.2	1.78×10^{-10}
6	SPATC1L	1.02×10^{-9}	AL645608.2	9.80×10^{-17}	SPATC1L	5.55×10^{-10}
7	RPL13A	3.62×10^{-9}	RNF219	6.23×10^{-13}	RPL13A	2.70×10^{-9}
8	ADRB3	1.07×10^{-8}	HLA-DRB1	1.38×10^{-12}	ADRB3	1.73×10^{-8}
9	TM4SF18	4.66×10^{-8}	OR7G3	3.01×10^{-12}	TM4SF18	3.41×10^{-8}
10	OR7G3	6.74×10^{-8}	EPS8L1	2.11×10^{-11}	OR7G3	7.35×10^{-8}

rank	conservation		log likelihoods		number of mutations	
	gene name	q-value	gene name	q-value	gene name	q-value
1	TP53	1.22×10^{-21}	TP53	1.18×10^{-21}	TP53	5.11×10^{-18}
2	STK11	2.24×10^{-3}	STK11	1.86×10^{-4}	STK11	9.48×10^{-3}
3	NF1	2.29×10^{-2}	NF1	1.20×10^{-2}	NF1	1.56×10^{-1}
4	CDKN2A	1.08×10^{-1}	PBRM1	6.48×10^{-2}	VHL	3.50×10^{-1}
5	VHL	1.53×10^{-1}	VHL	6.48×10^{-2}	PBRM1	3.50×10^{-1}
6	PBRM1	1.91×10^{-1}	ITGA11	7.60×10^{-2}	XCL1	5.12×10^{-1}
7	XCL1	2.54×10^{-1}	SKA3	1.25×10^{-1}	ITGA11	6.88×10^{-1}
8	SKA3	2.54×10^{-1}	RP13-512J5.1	1.38×10^{-1}	B2M	1.00
9	ITAG11	2.80×10^{-1}	PTEN	1.66×10^{-1}	MARCH1	1.00
10	B2M	2.80×10^{-1}	XCL1	1.78×10^{-1}	TGFB1	1.00

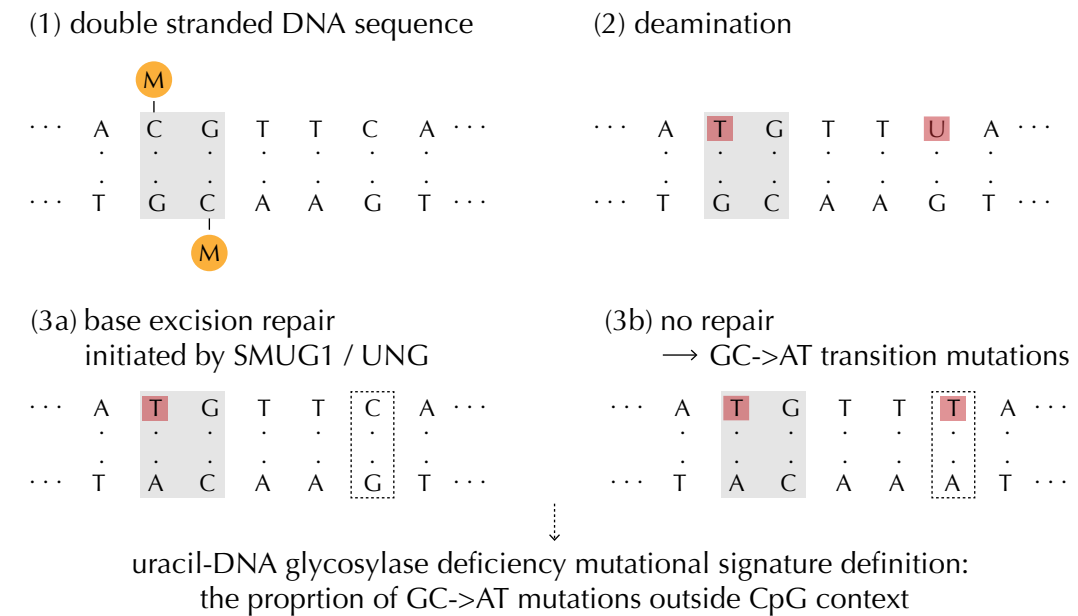
gene name	TERT	FRG2B	HLA-DRB5	MUC3A	AL645608.2	SPATC1L	RPL13A	ADRB3	TM4SF18	OR7G3	HLA-DRB1	RNF219	EPS8L1
region size (bp)	600	754	835	500	982	894	907	500	749	500	792	867	2339
observed mutations	77	74	43	17	41	26	27	18	19	20	33	32	55
expected mutations	3.4	3.9	4.1	0.9	4.0	3.1	3.0	1.7	2.1	2.2	4.0	5.9	6.8

gene name	TP53	STK11	NF1	CDKN2A	VHL	PBRM1	XCL1	SKA3	ITGA11	B2M	MARCH1	TGFB1	PTEN	RP13-512J5.1
region size (bp)	52	36	232	24	8	124	8	32	116	10	36	24	32	6
observed mutations	12	4	5	2	2	4	2	2	4	2	3	2	2	1
expected mutations	0.1	0.1	0.3	0.0	0.0	0.2	0.0	0.1	0.3	0.0	0.2	0.0	0.1	0.0

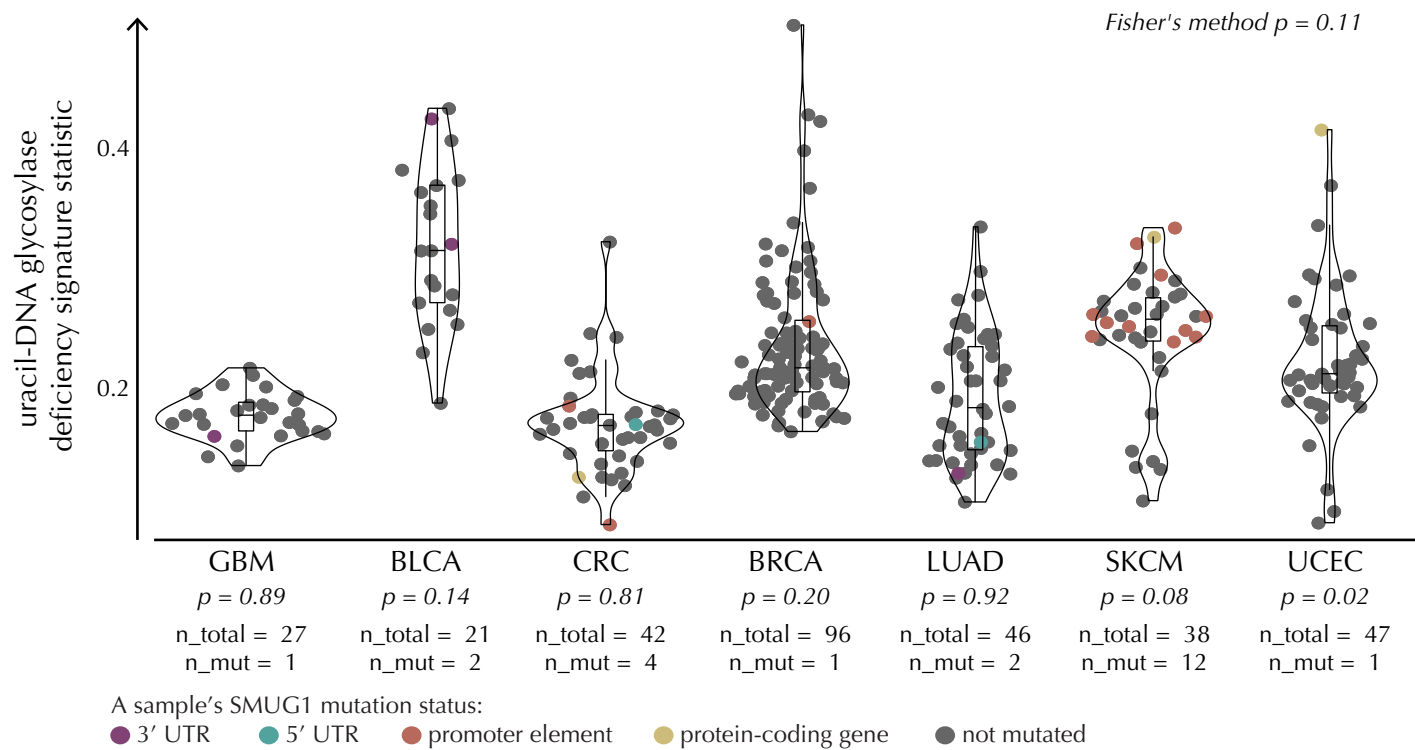
A



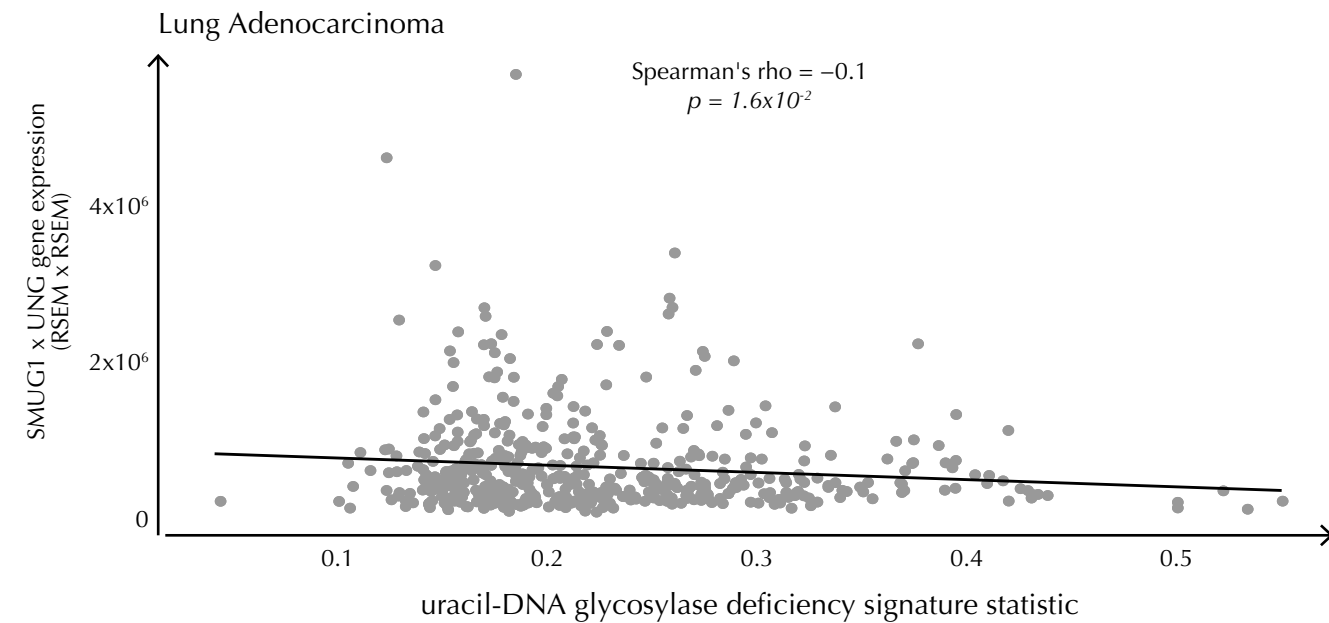
B

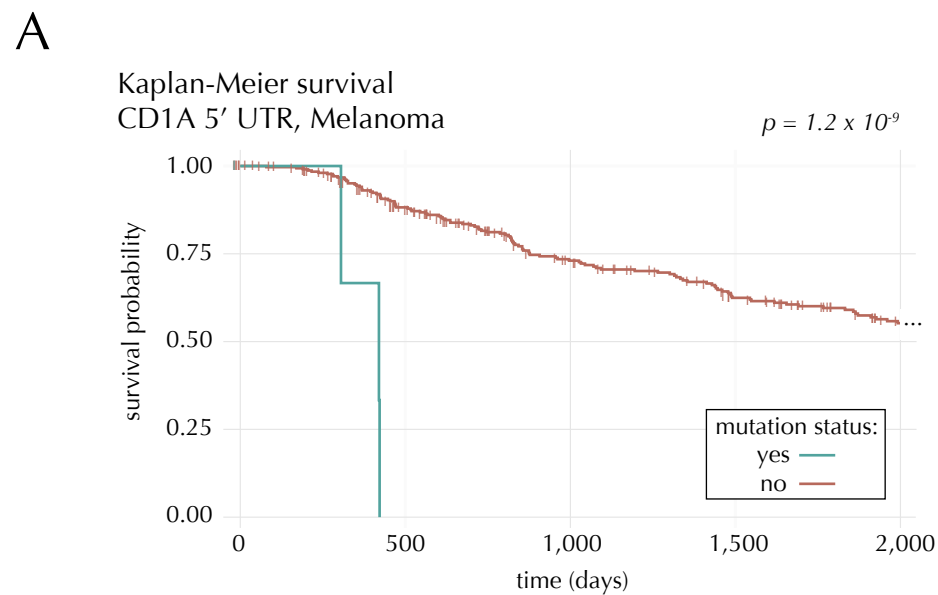


C



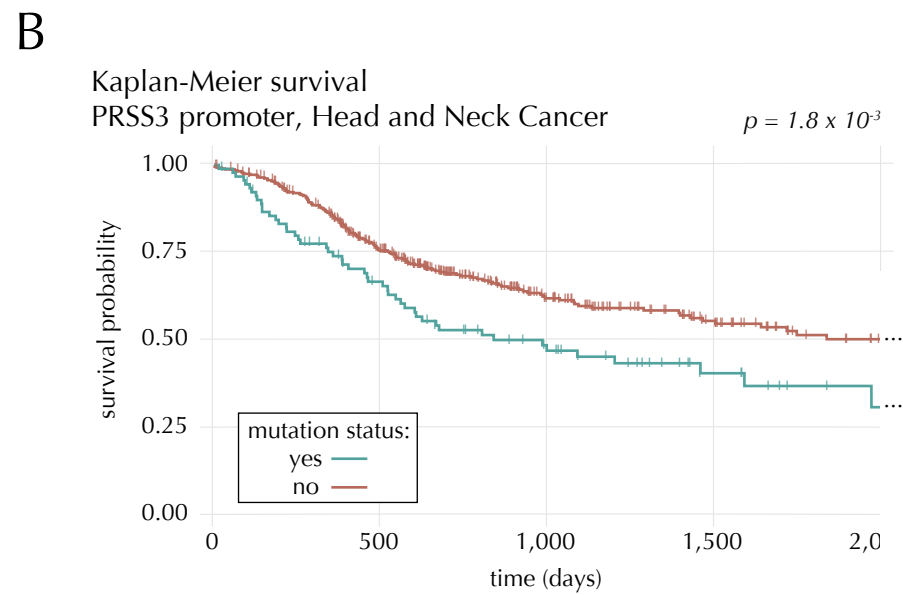
D





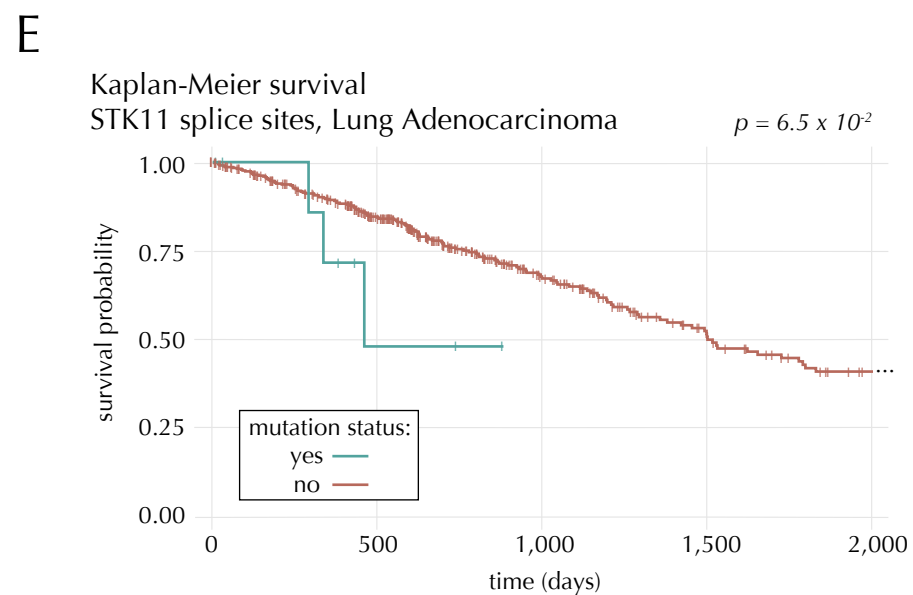
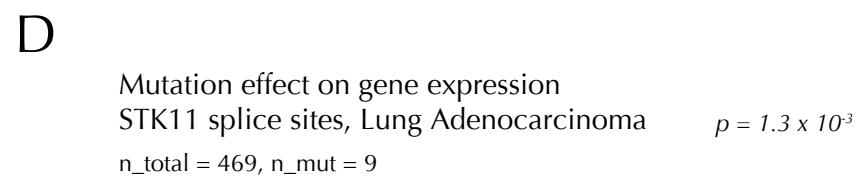
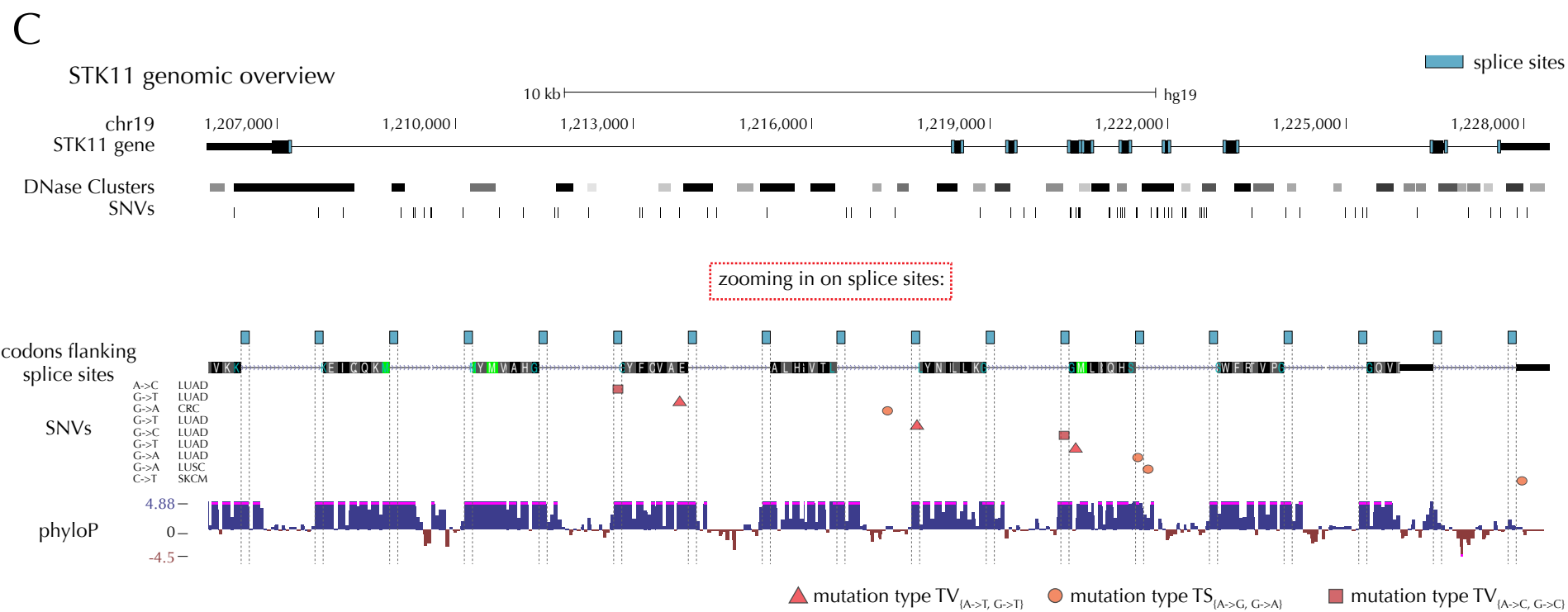
Number of patients at risk

follow-up time (days)	0	500	1,000	1,500	2,000
mutation status: no	320	250	177	136	98
mutation status: yes	4	0	0	0	0



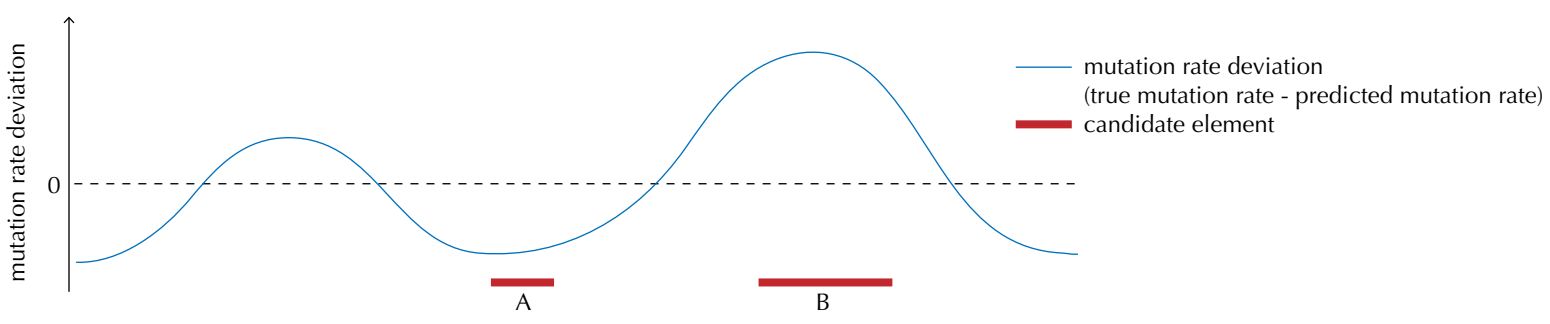
Number of patients at risk

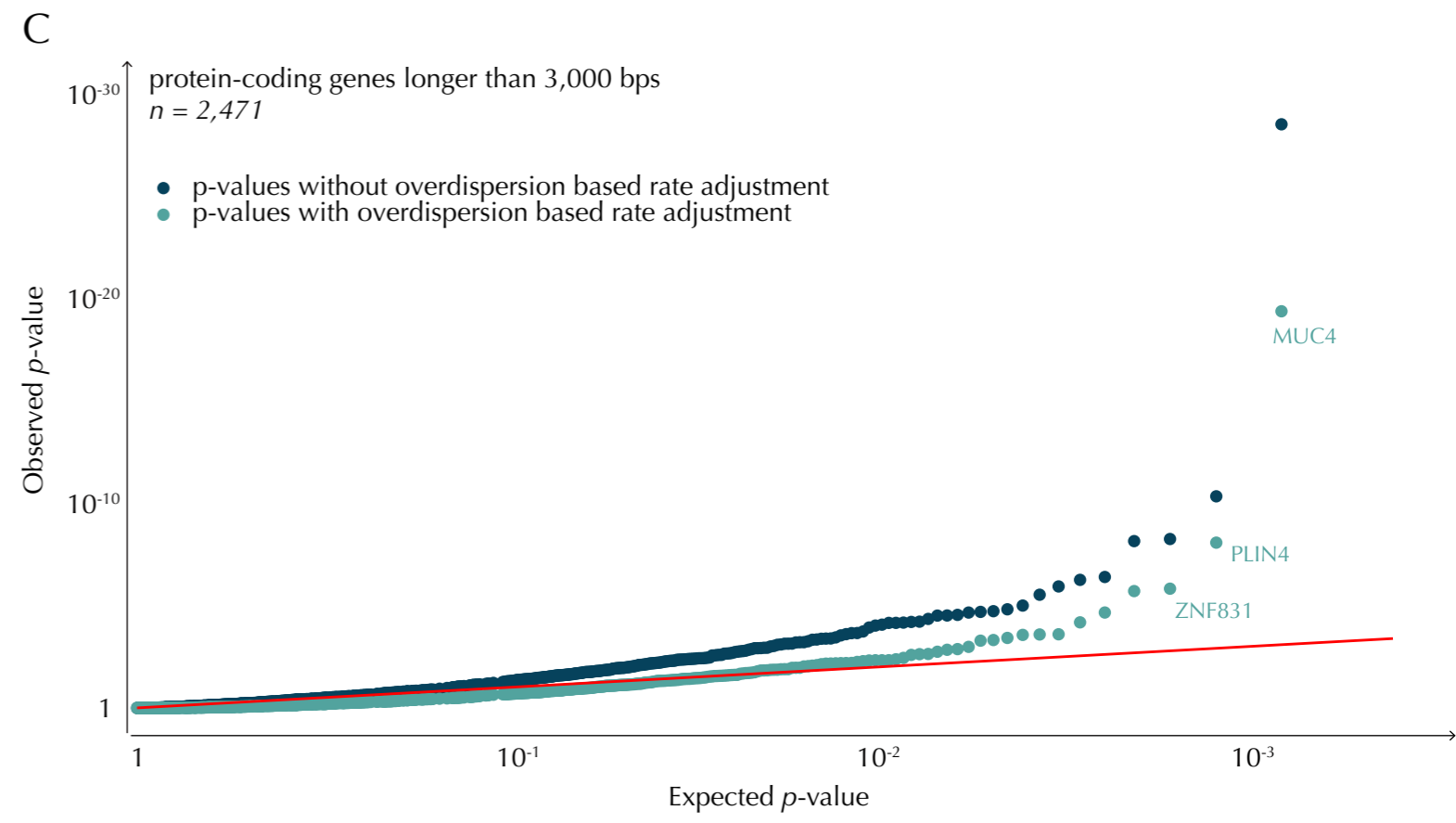
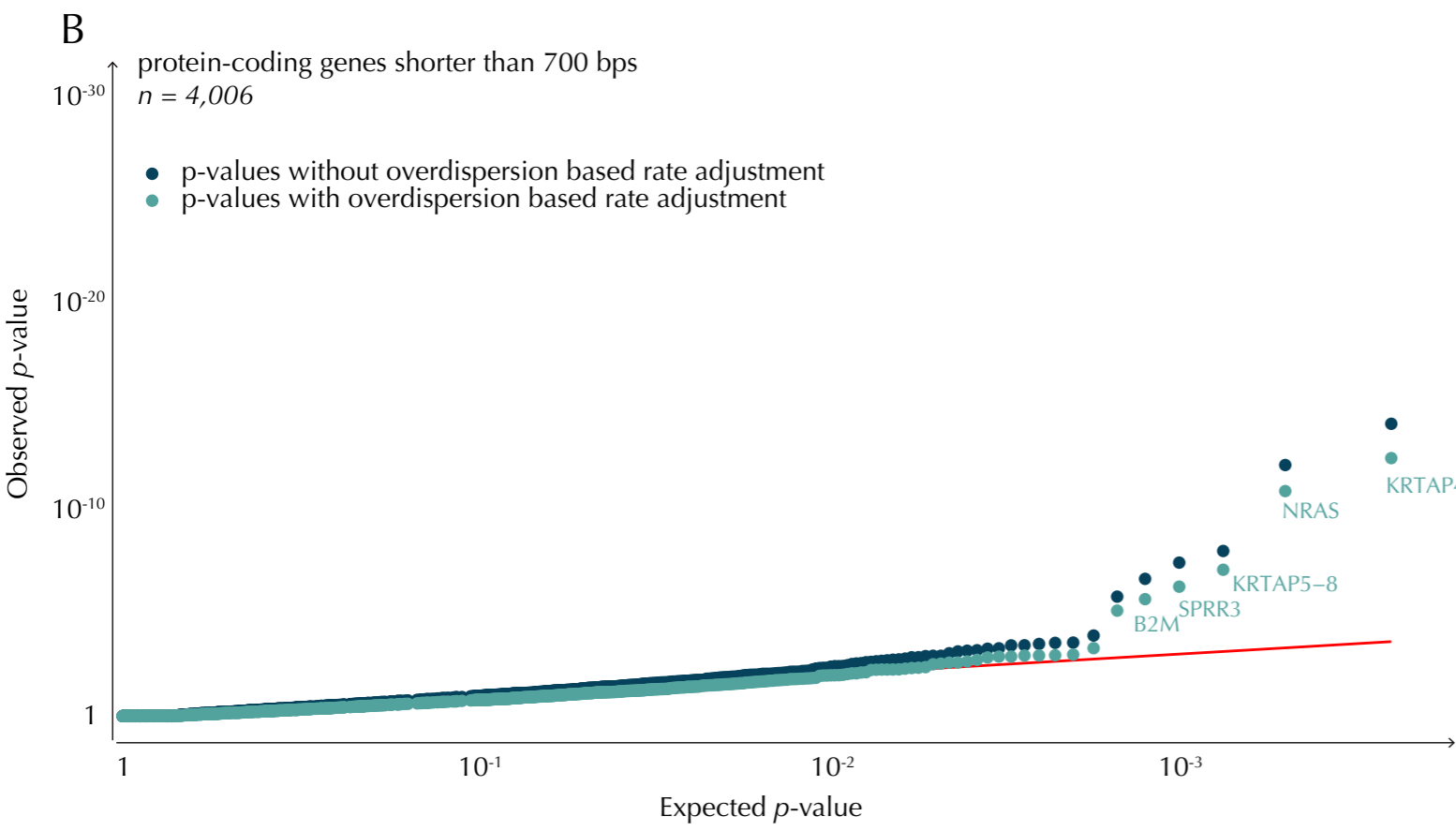
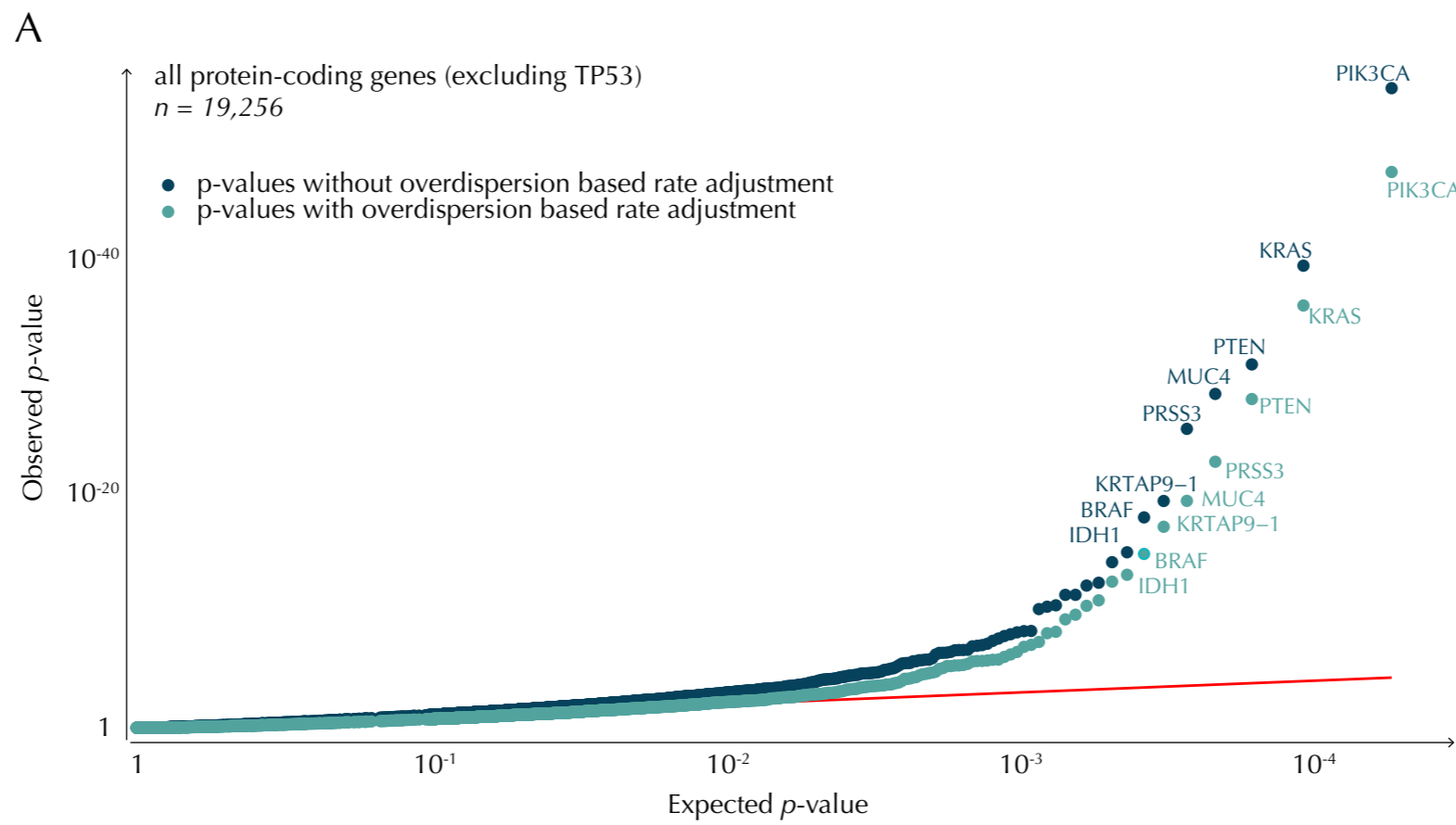
follow-up time (days)	0	500	1,000	1,500	2,000
mutation status: no	393	251	121	67	37
mutation status: yes	91	53	30	13	5

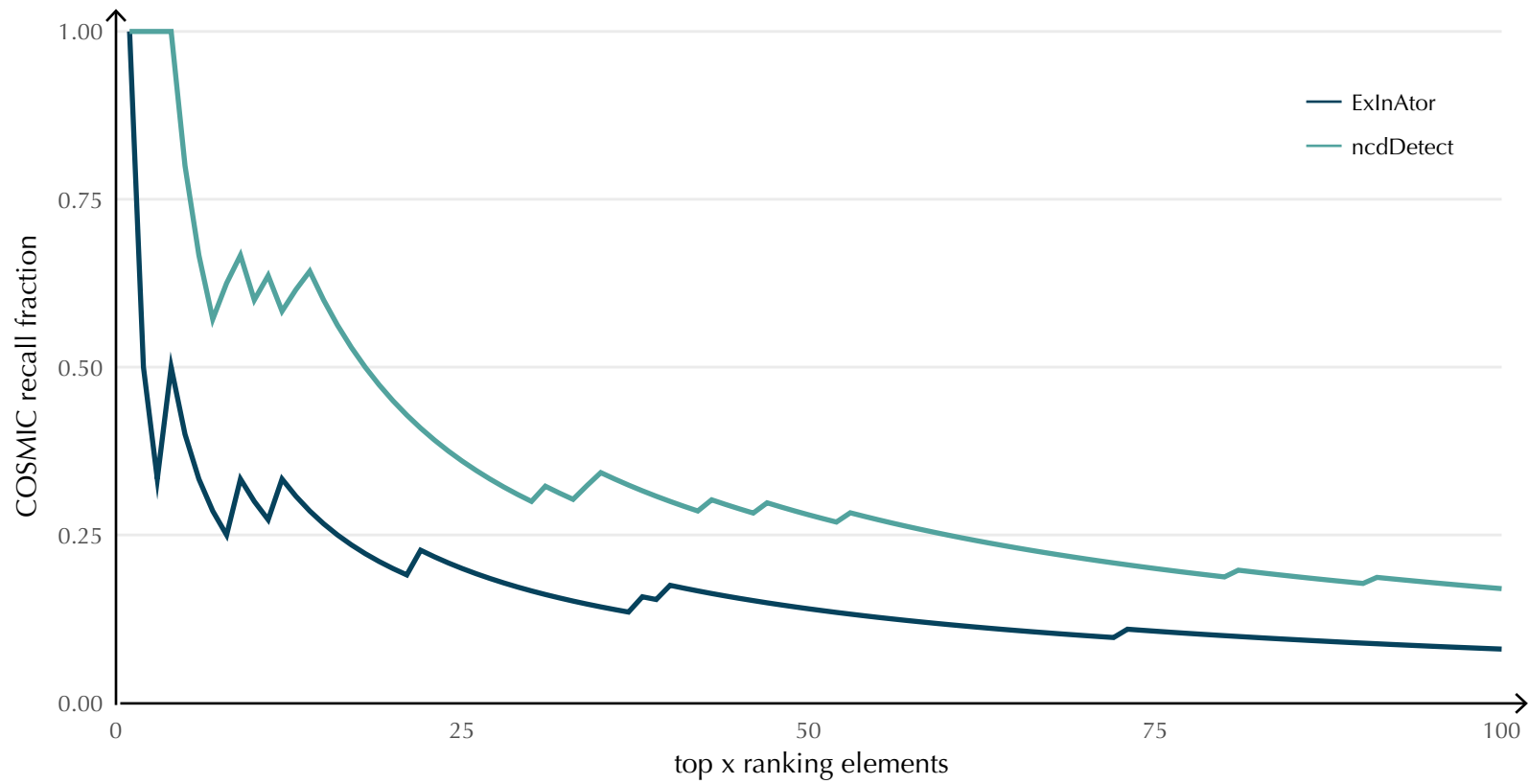


Number of patients at risk

follow-up time (days)	0	500	1,000	1,500	2,000
mutation status: no	429	281	121	61	36
mutation status: yes	9	2	0	0	0

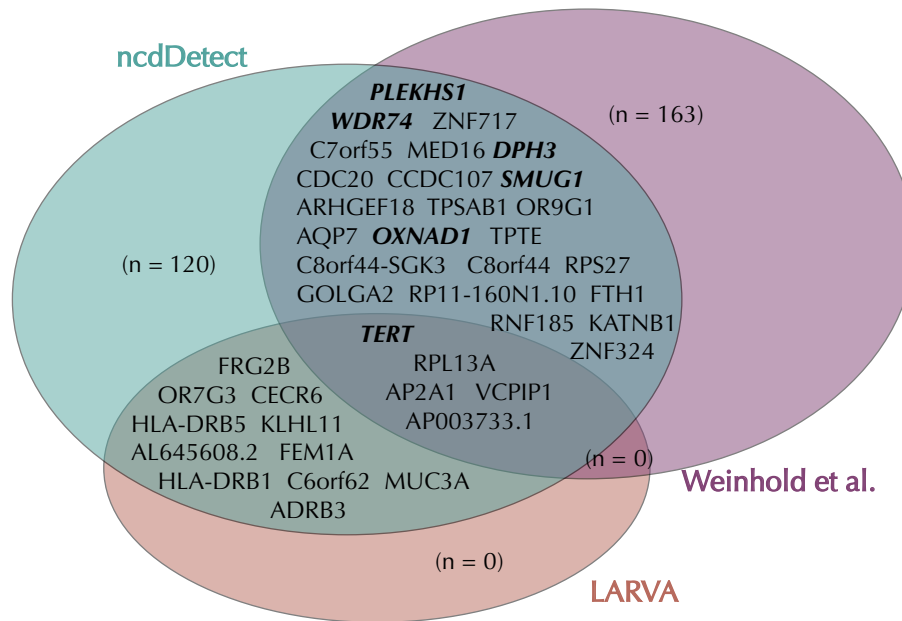






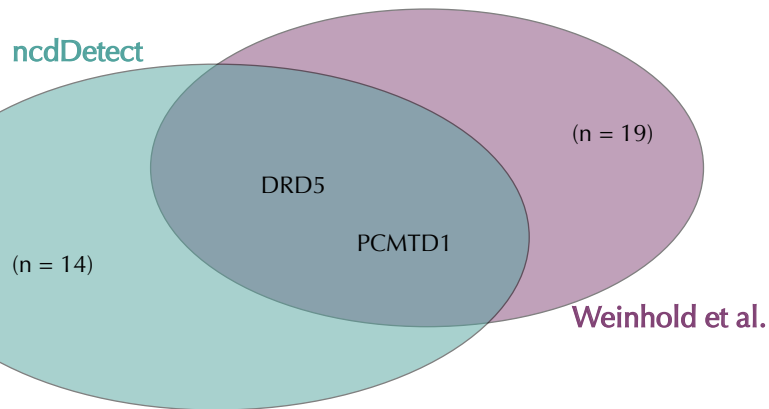
A

promoter elements



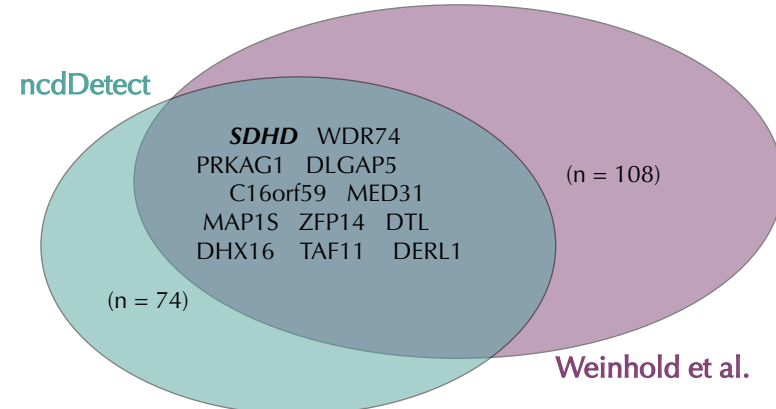
B

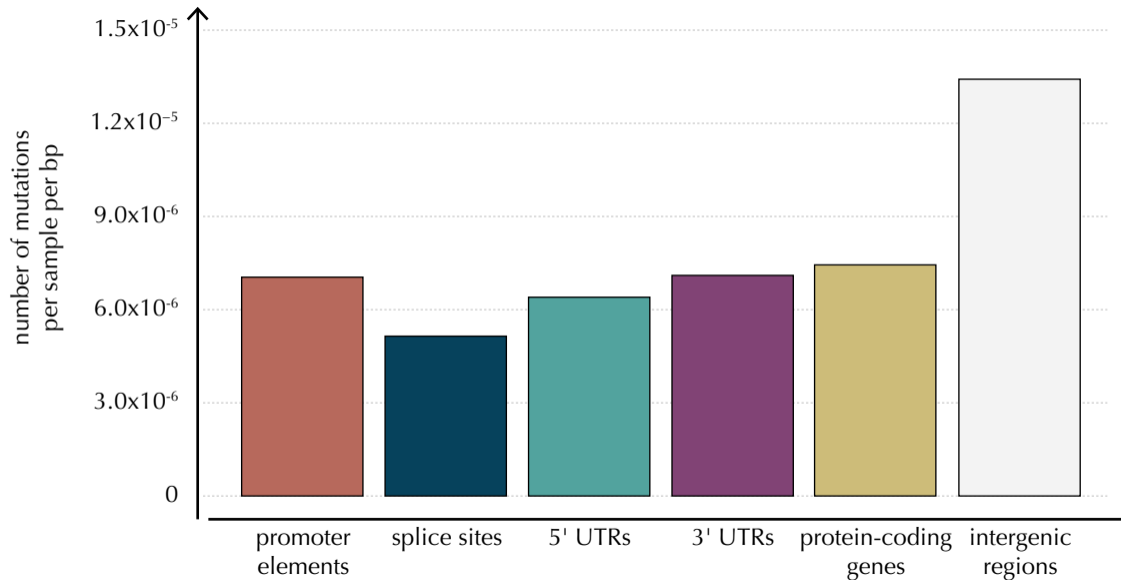
3' UTRs



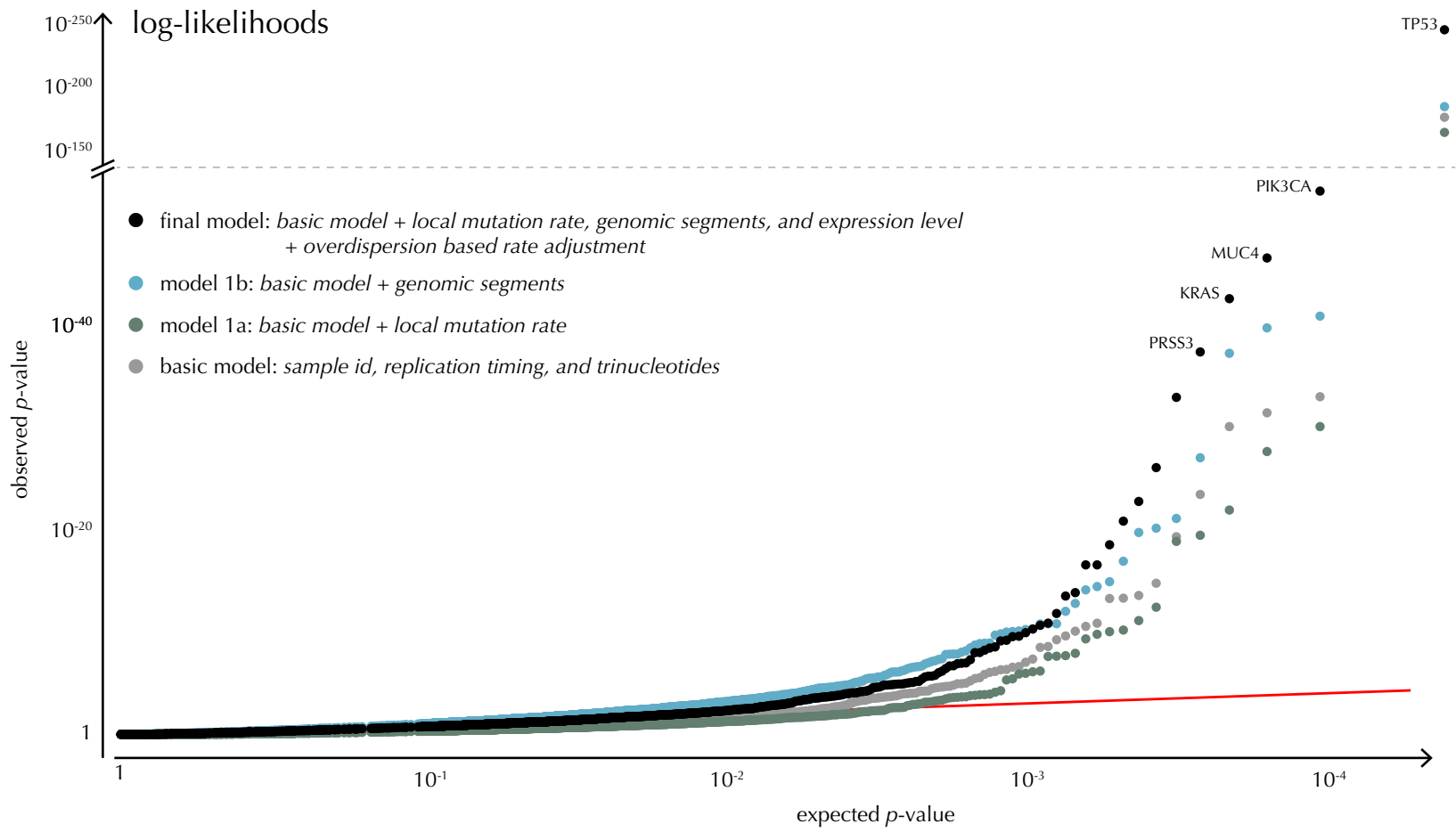
C

5' UTRs

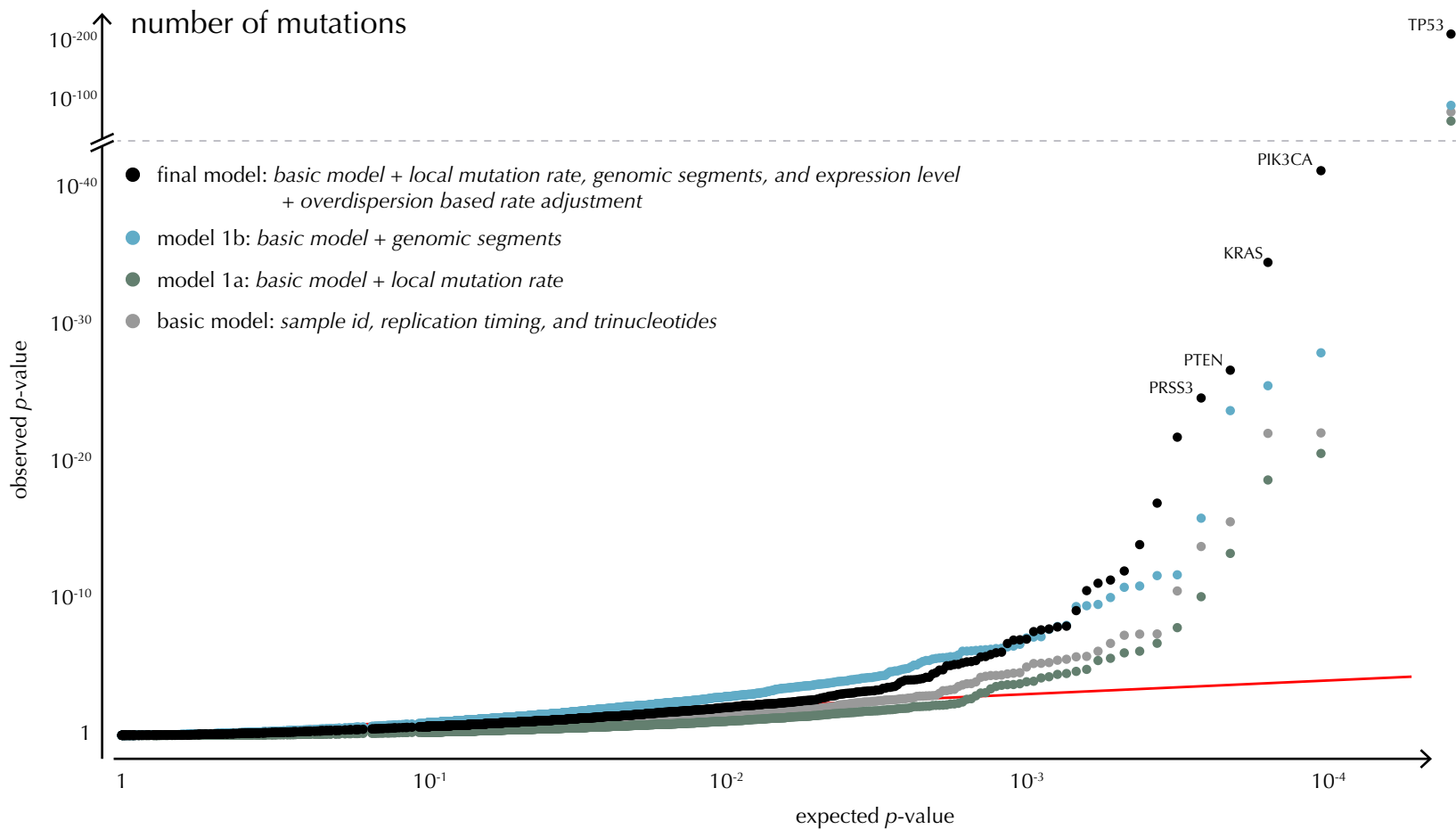


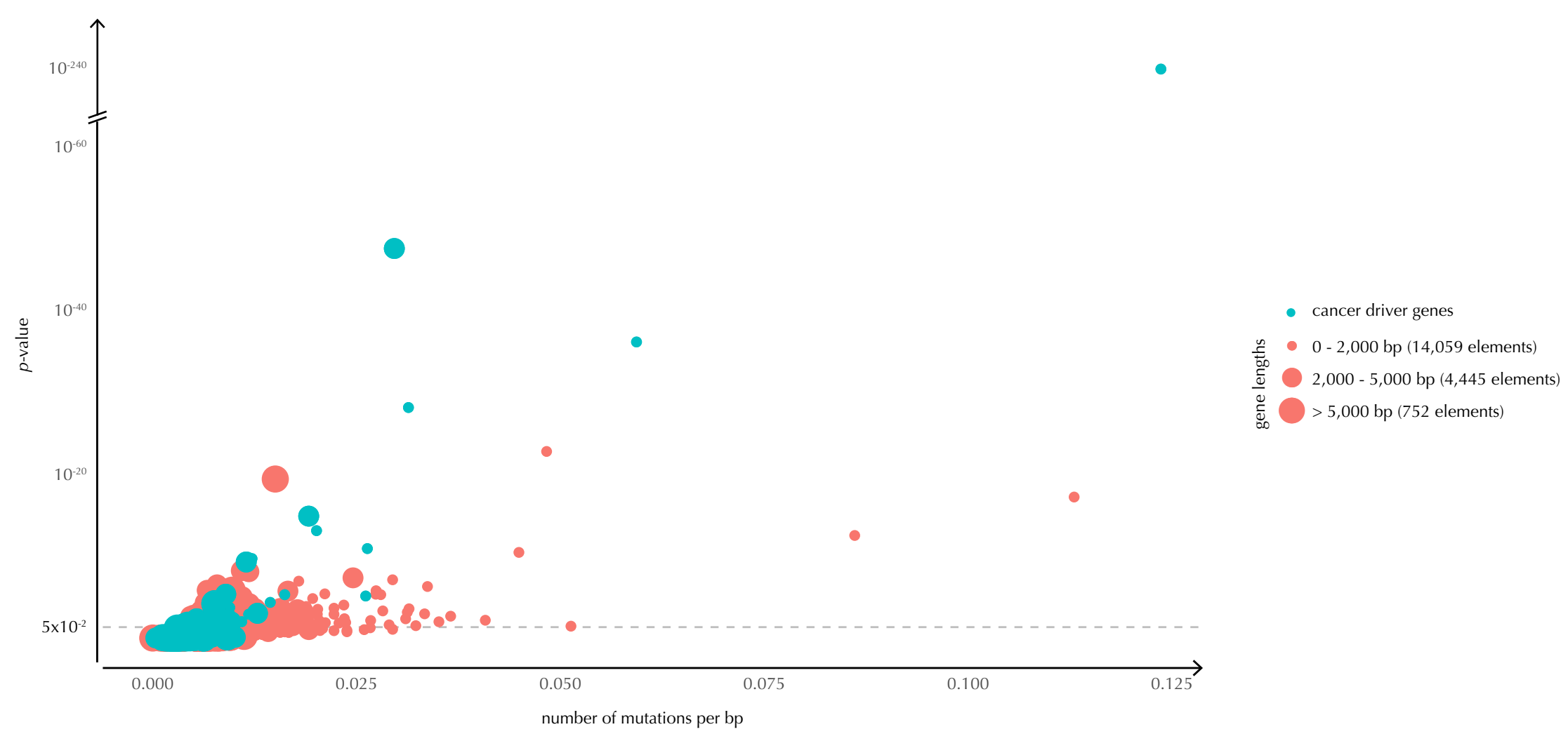


A

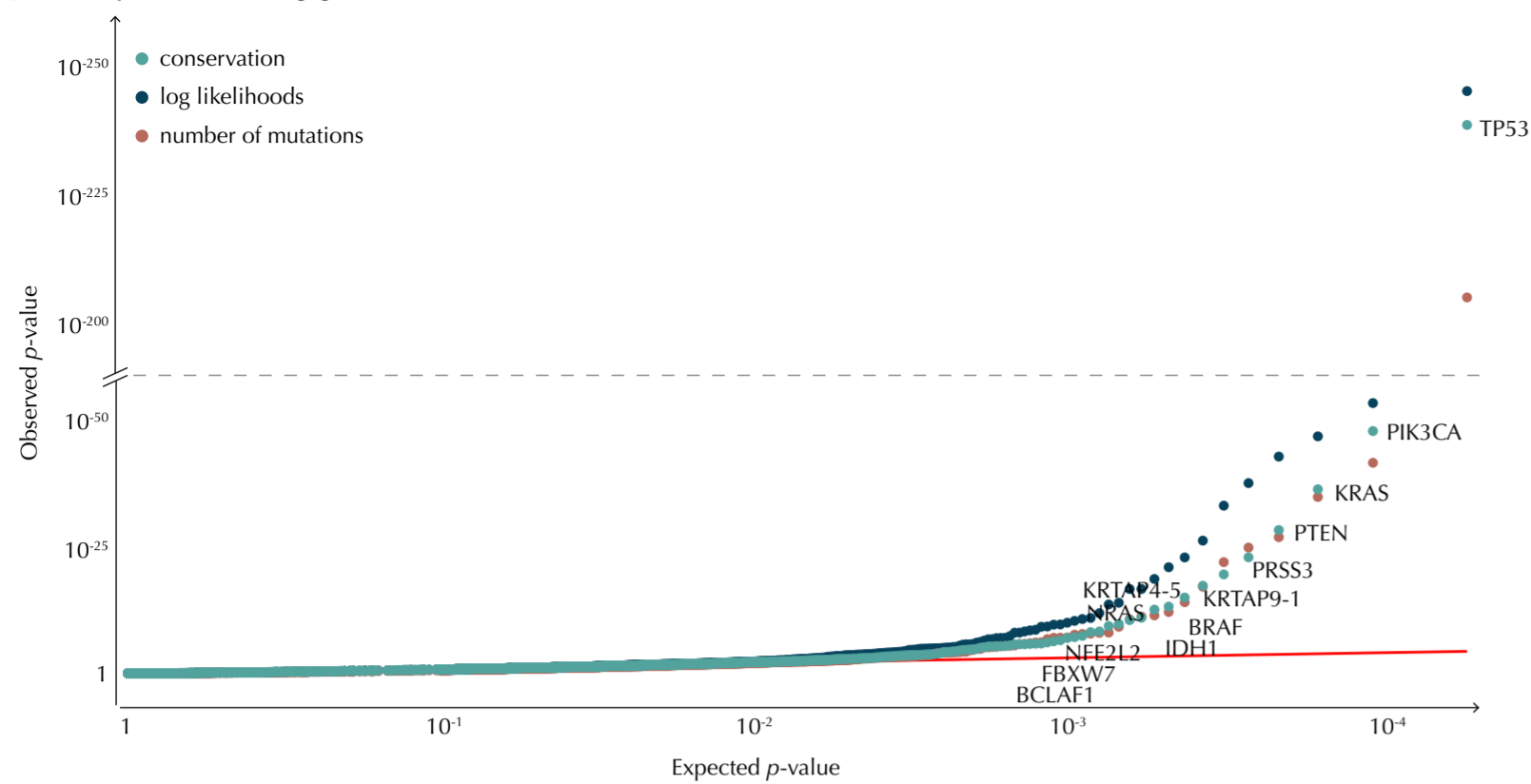


B





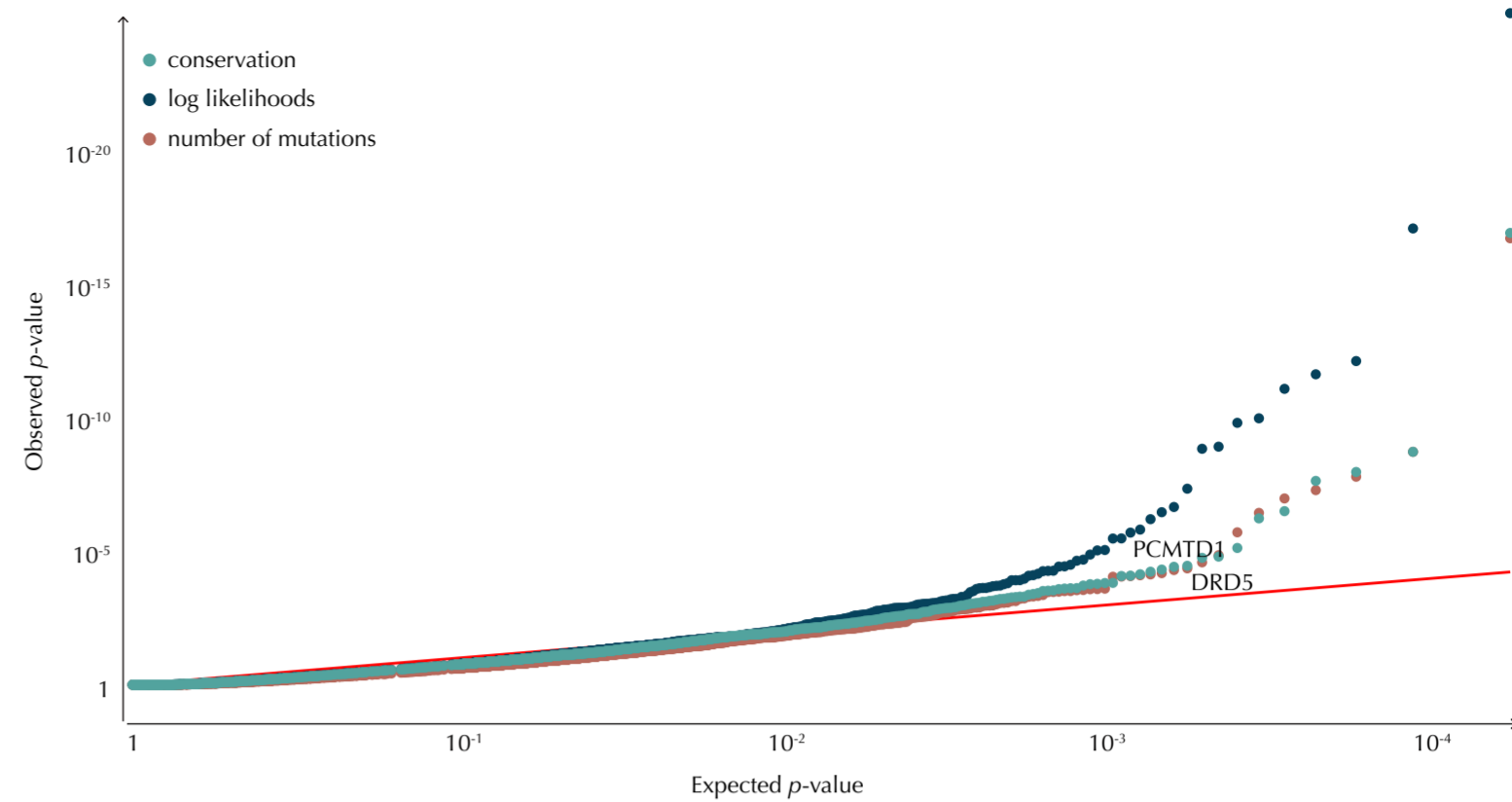
A protein-coding genes ($n = 19,256$)



rank	conservation		log likelihoods		number of mutations	
	gene name	q-value	gene name	q-value	gene name	q-value
1	TP53	1.15×10^{-235}	TP53	2.65×10^{-242}	TP53	8.21×10^{-202}
2	PIK3CA	2.82×10^{-44}	PIK3CA	8.90×10^{-50}	PIK3CA	4.69×10^{-38}
3	KRAS	4.93×10^{-33}	MUC4	2.01×10^{-43}	KRAS	1.61×10^{-31}
4	PTEN	3.70×10^{-25}	KRAS	1.42×10^{-39}	PTEN	9.30×10^{-24}
5	PRSS3	6.80×10^{-20}	PRSS3	1.79×10^{-34}	PRSS3	8.07×10^{-22}
6	MUC4	1.26×10^{-16}	PTEN	4.07×10^{-30}	MUC4	4.95×10^{-19}
7	KRTAP9-1	1.77×10^{-14}	KRTAP9-1	2.51×10^{-23}	KRTAP9-1	2.83×10^{-14}
8	BRAF	3.25×10^{-12}	KRTAP4-5	4.42×10^{-20}	KRTAP4-5	2.71×10^{-11}
9	IDH1	1.76×10^{-10}	BRAF	3.28×10^{-18}	BRAF	2.04×10^{-9}
10	KRTAP4-5	6.11×10^{-10}	AL390778.1	5.96×10^{-16}	IDH1	8.41×10^{-9}

gene name	TP53	PIK3CA	KRAS	PTEN	PRSS3	MUC4	KRTAP9-1	BRAF	IDH1	KRTAP4-5	AL390778.1
region size (bp)	1,378	3,207	708	1,212	1,056	16,239	770	2,301	1,245	546	735
observed mutations	169	95	42	38	51	244	87	44	25	47	33
expected mutations	3.3	12.0	2.2	2.5	3.5	74.1	2.8	9.5	3.5	2.4	5.2

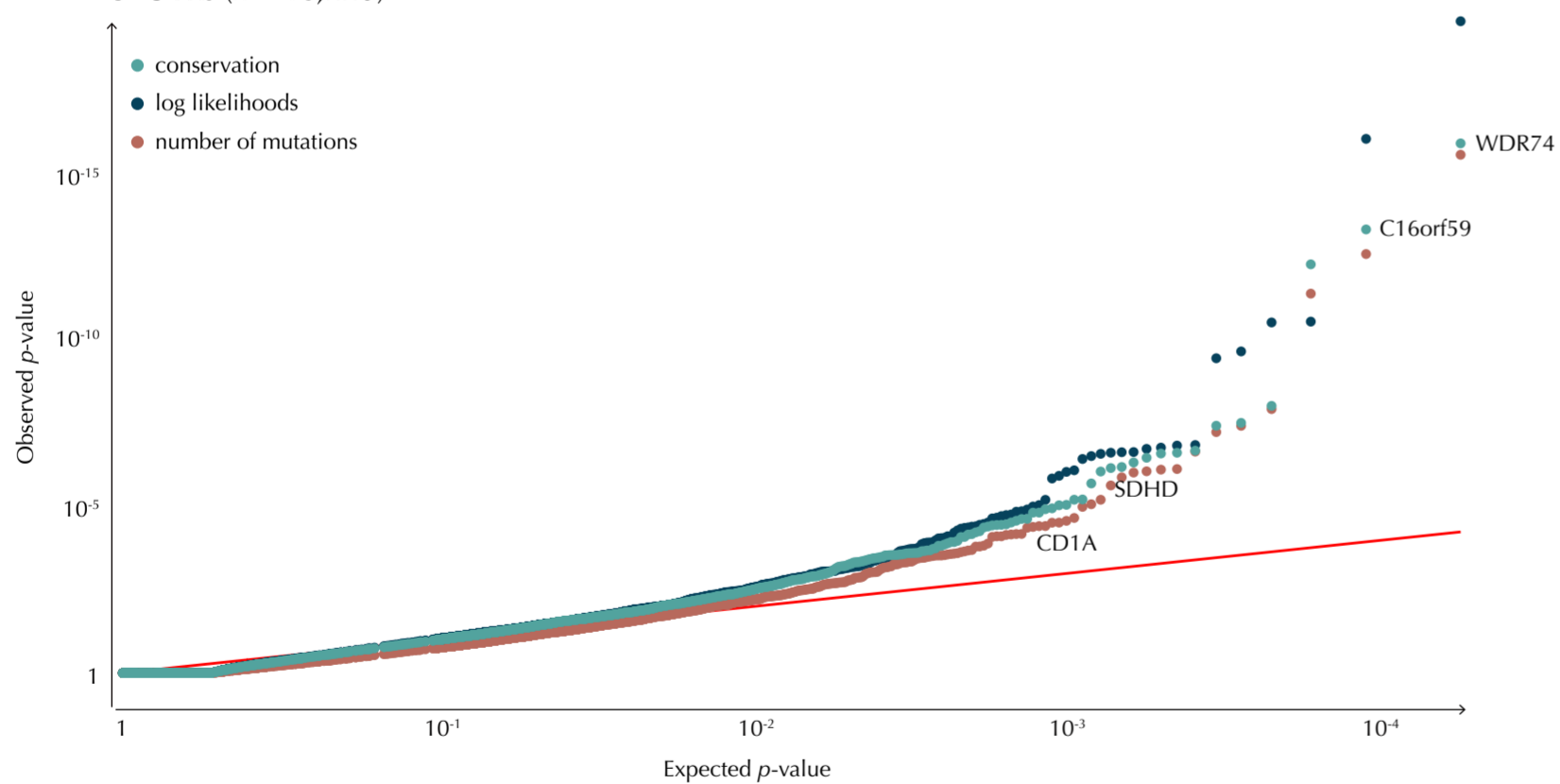
B 3' UTRs ($n = 18,481$)



rank	conservation		log likelihoods		number of mutations	
	gene name	q-value	gene name	q-value	gene name	q-value
1	MYO5B	1.46×10^{-13}	MYO5B	7.03×10^{-22}	MYO5B	2.33×10^{-13}
2	PRSS3	1.41×10^{-5}	VPS53	4.95×10^{-14}	PRSS3	1.41×10^{-5}
3	ADD2	5.44×10^{-5}	PRSS3	3.44×10^{-9}	ADD2	8.21×10^{-5}
4	SEC14L1	8.87×10^{-5}	SEC14L1	8.18×10^{-9}	SEC14L1	1.98×10^{-4}
5	FAHD2B	9.92×10^{-4}	TBC1D22A	2.33×10^{-8}	FAHD2B	3.27×10^{-4}
6	VPS53	1.55×10^{-3}	FAHD2B	2.53×10^{-7}	VPS53	9.61×10^{-4}
7	LRTM1	1.75×10^{-2}	ADD2	3.21×10^{-7}	TBC1D22A	4.48×10^{-3}
8	SGCZ	3.23×10^{-2}	FAM101B	2.24×10^{-6}	LRTM1	2.89×10^{-2}
9	ELOVL3	3.23×10^{-2}	MUC19	2.40×10^{-6}	SGCZ	4.77×10^{-2}
10	DRD5	5.76×10^{-2}	LRTM1	7.01×10^{-5}	C7	7.25×10^{-2}

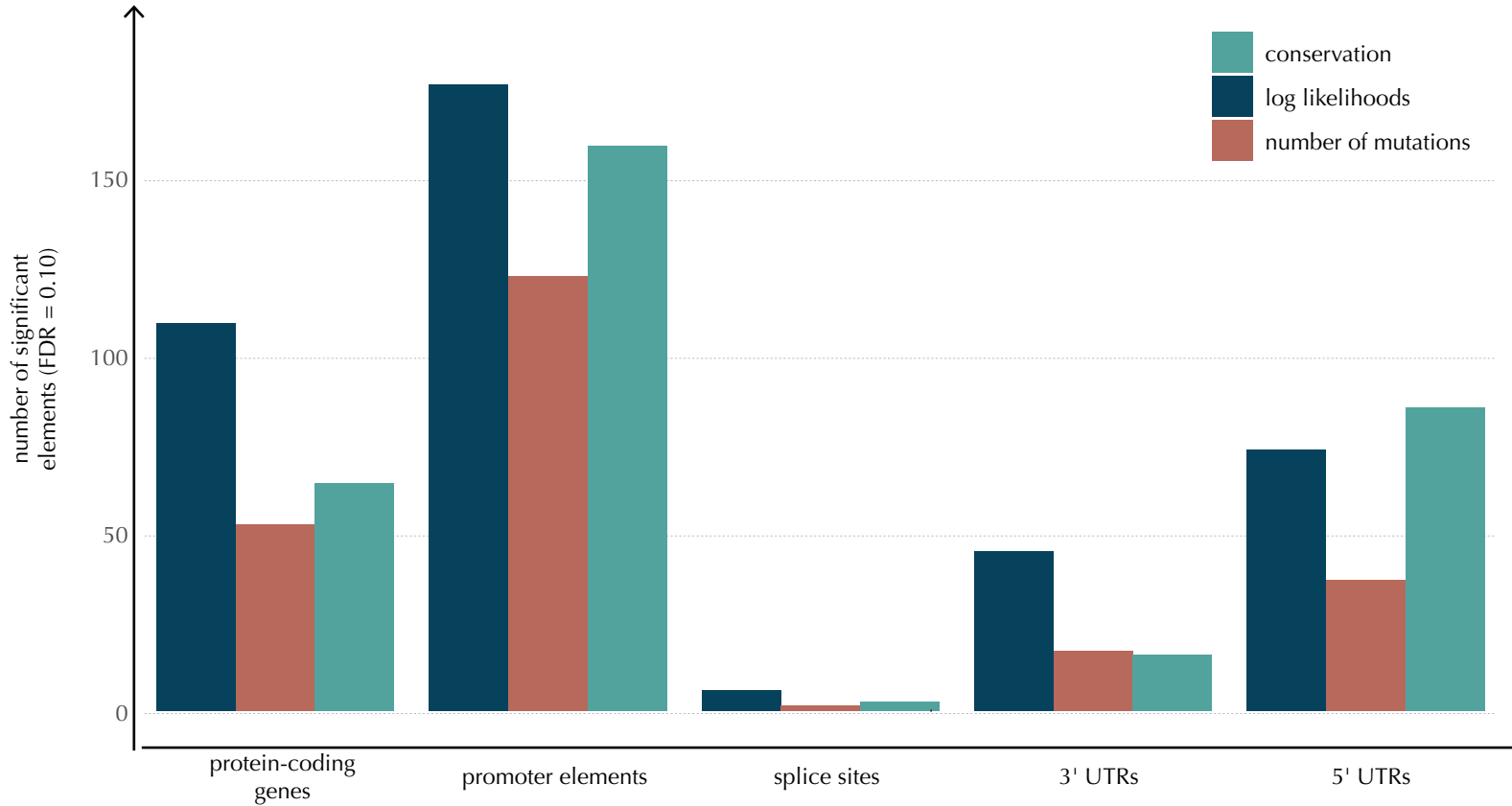
gene name	MYO5B	PRSS3	ADD2	SEC14L1	FAHD2B	VPS53	LRTM1	SGCZ	ELOVL3	DRD5	C7	TBC1D22A	FAM101B	MUC19
region size (bp)	3,658	52	8,257	3,052	172	11,085	222	579	337	545	1,366	2,067	3,186	8,002
observed mutations	78	7	50	31	10	56	9	19	8	22	22	37	19	126
expected mutations	10.1	0.2	12.4	7.9	0.8	19.0	1.1	5.0	1.0	3.5	6.1	9.9	4.6	61.6

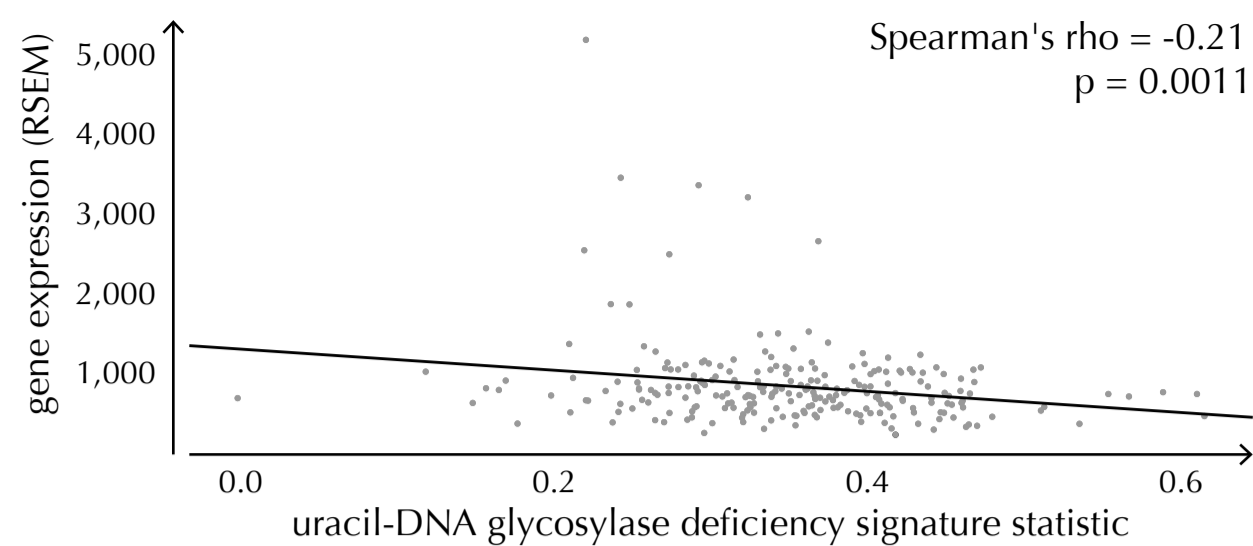
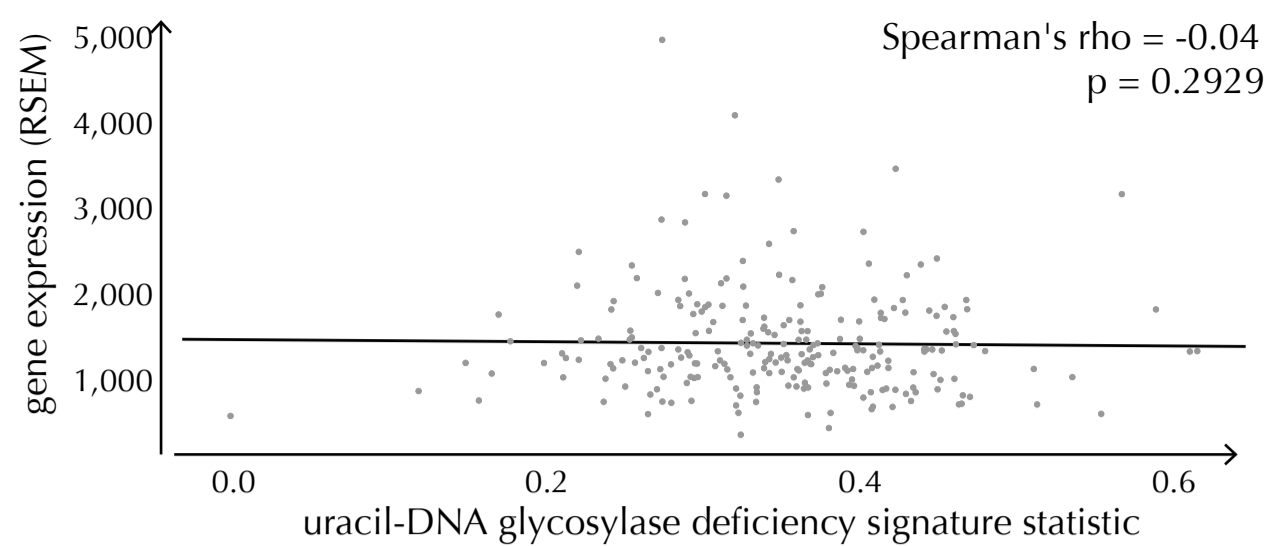
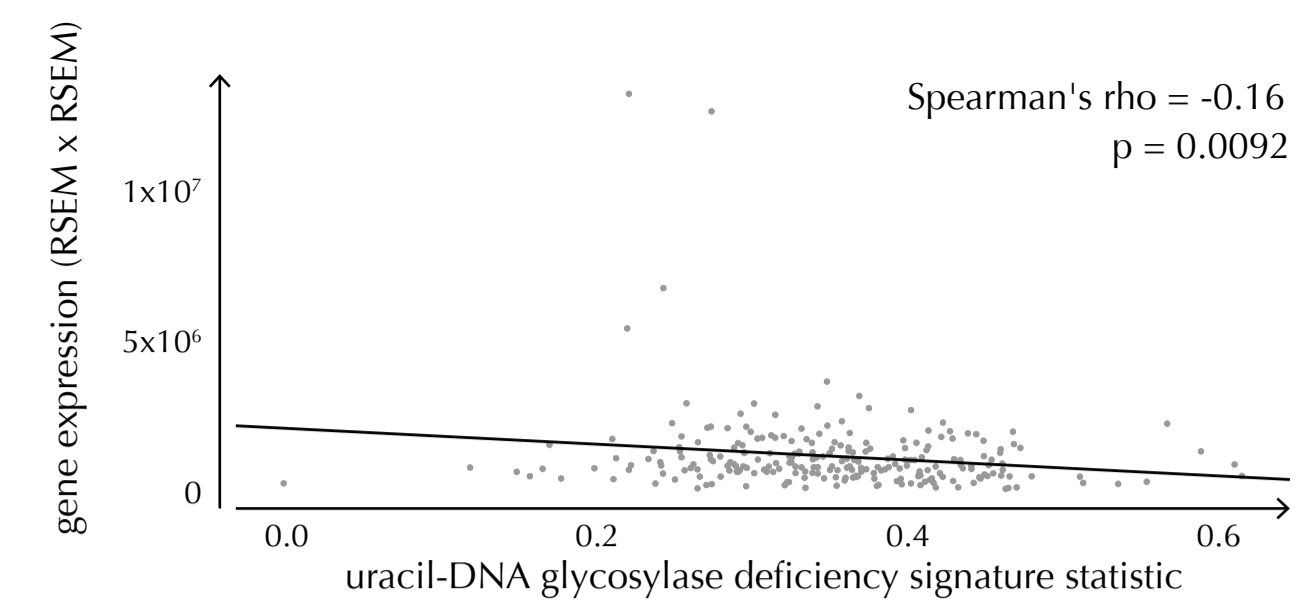
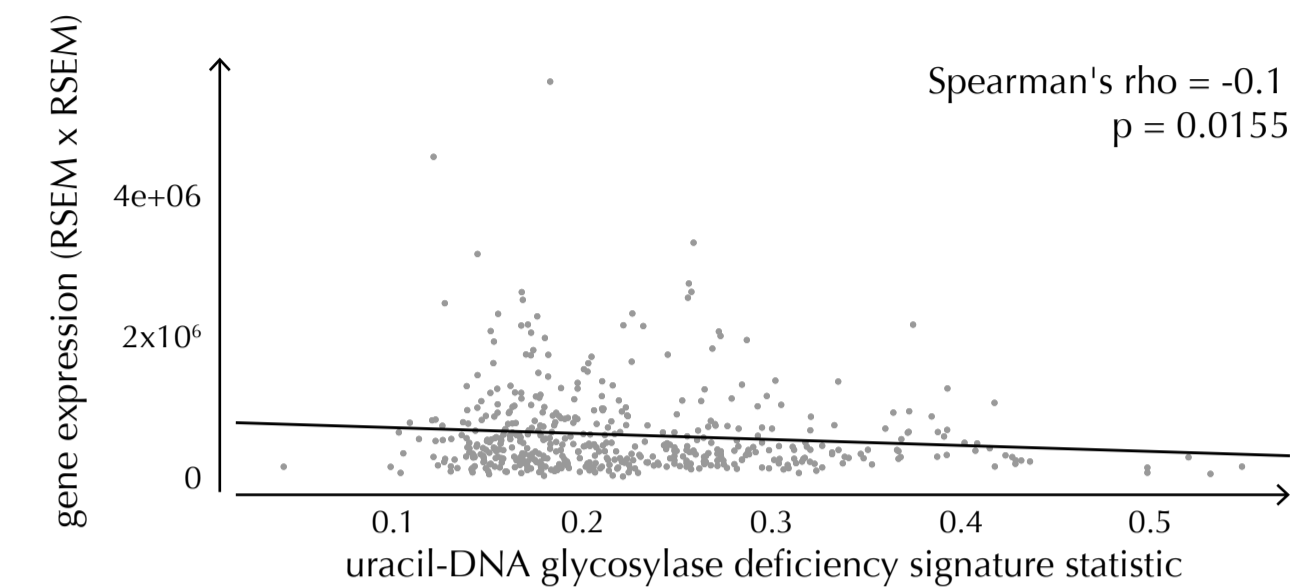
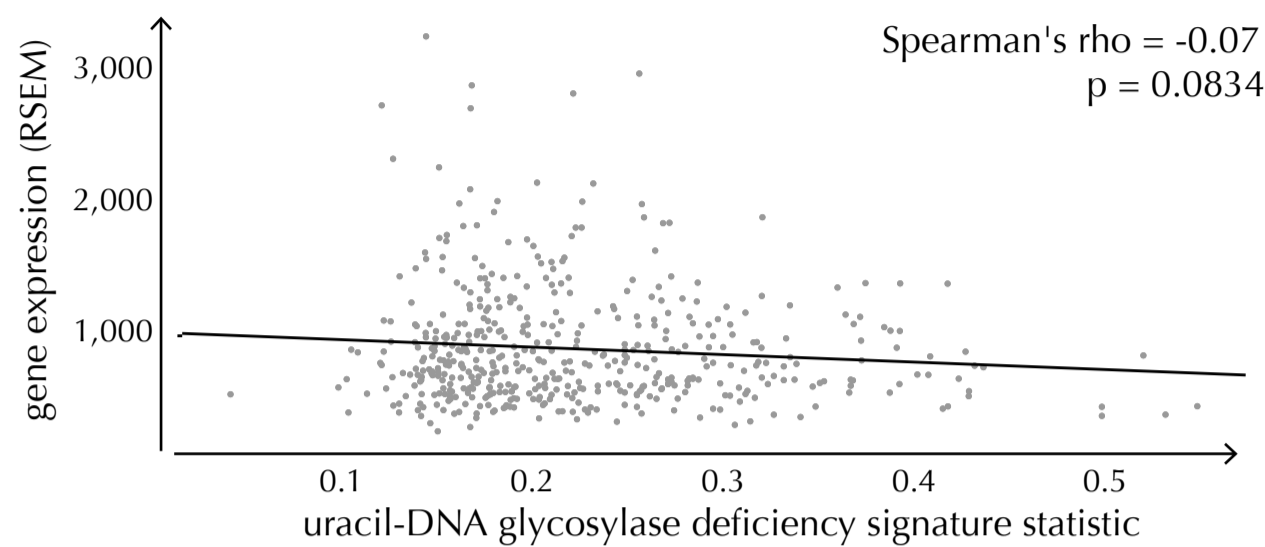
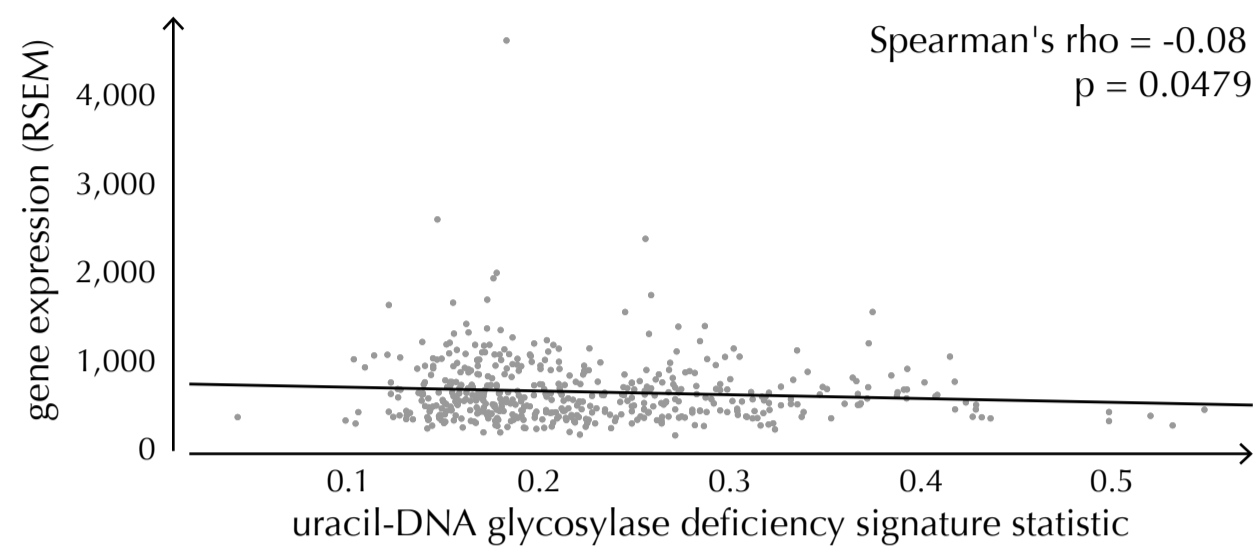
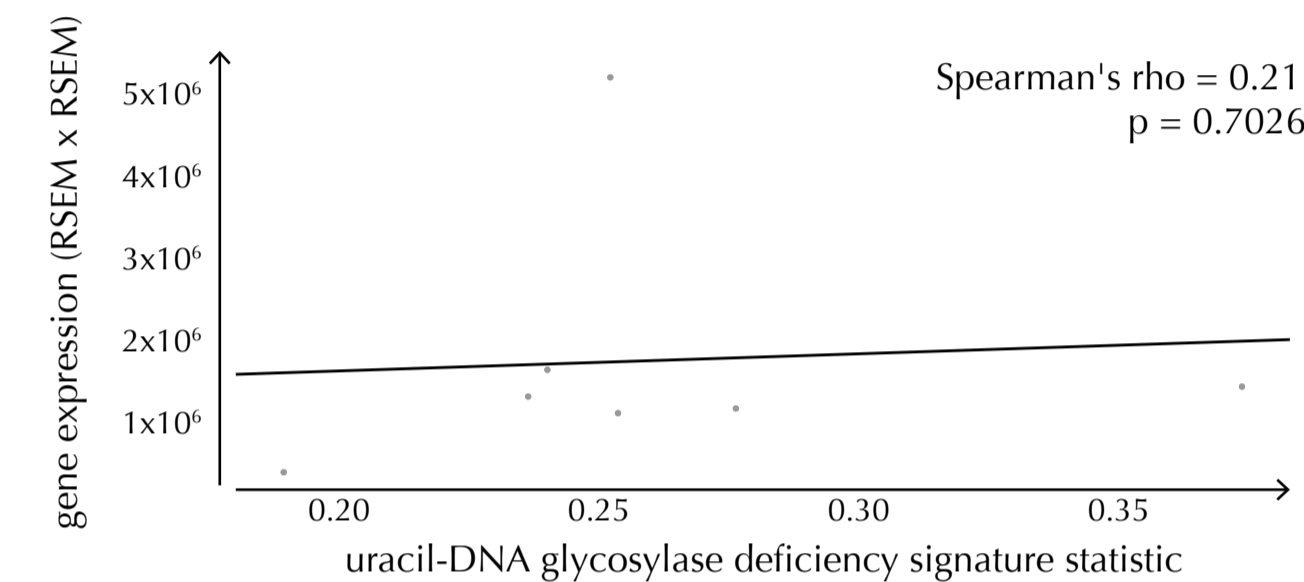
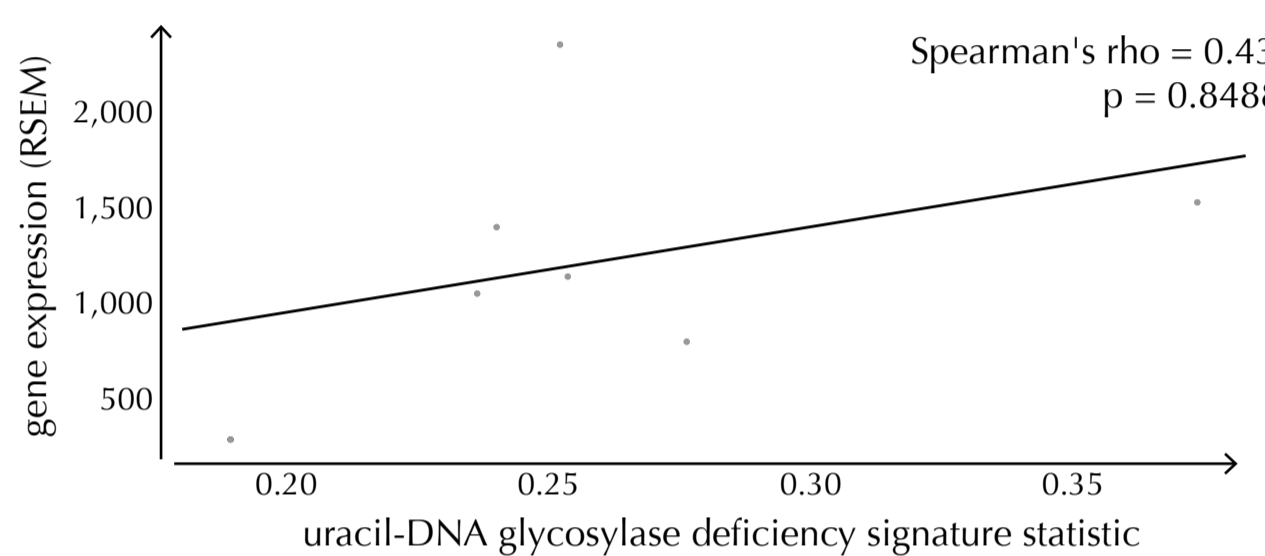
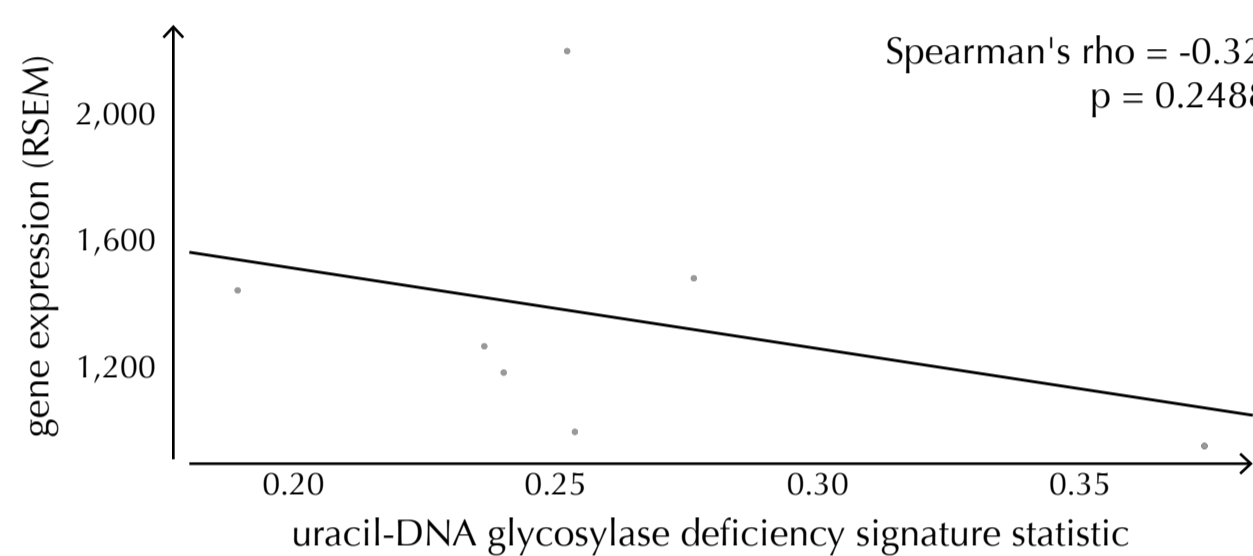
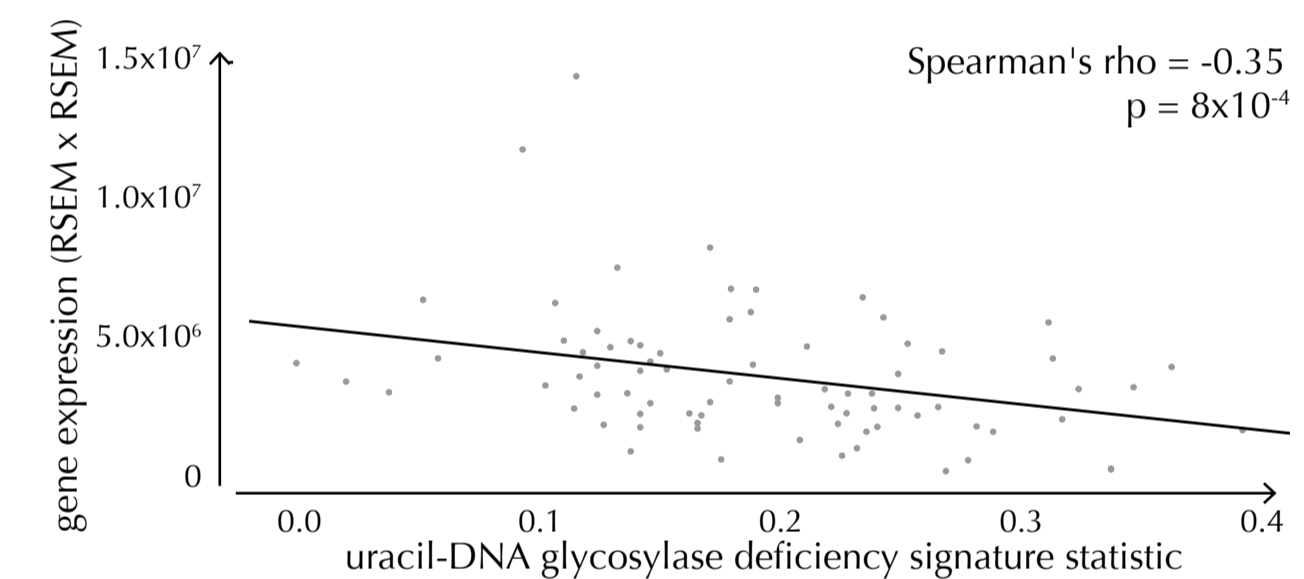
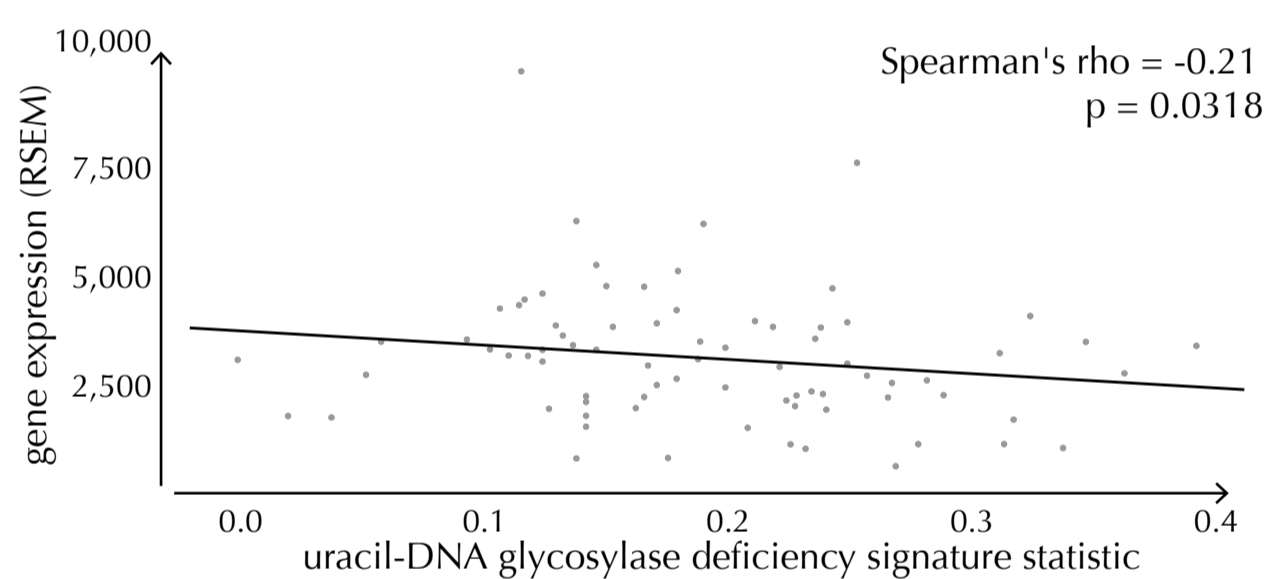
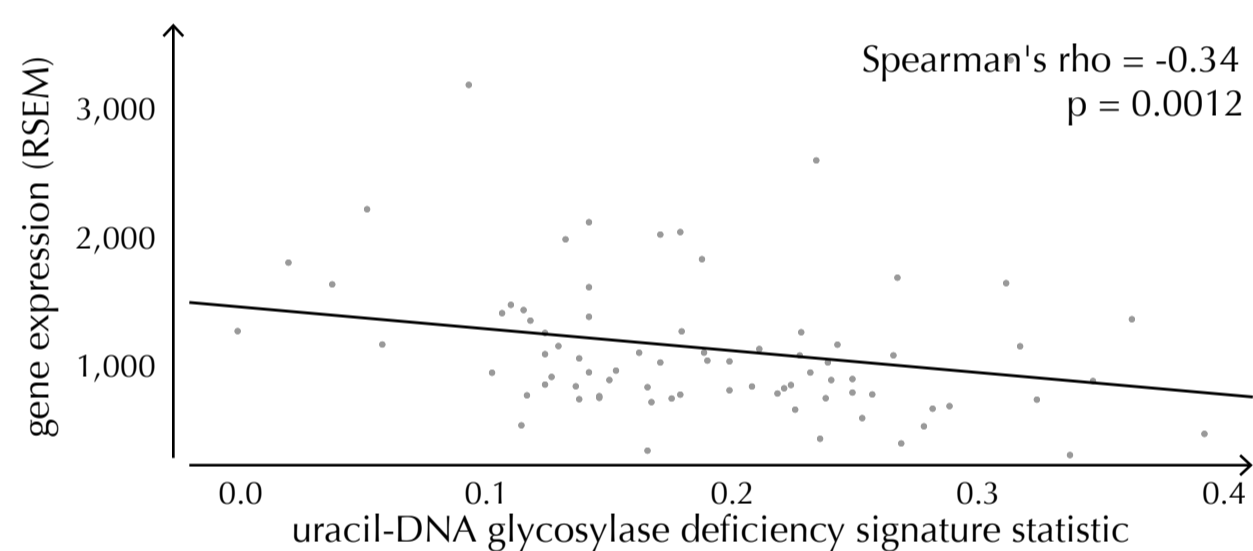
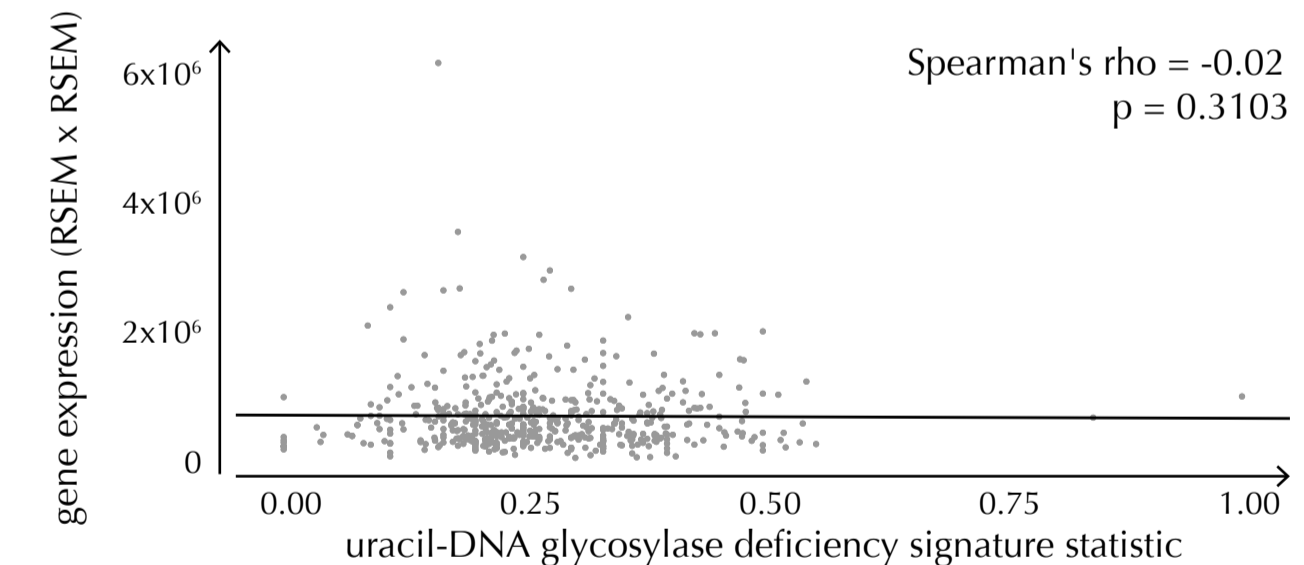
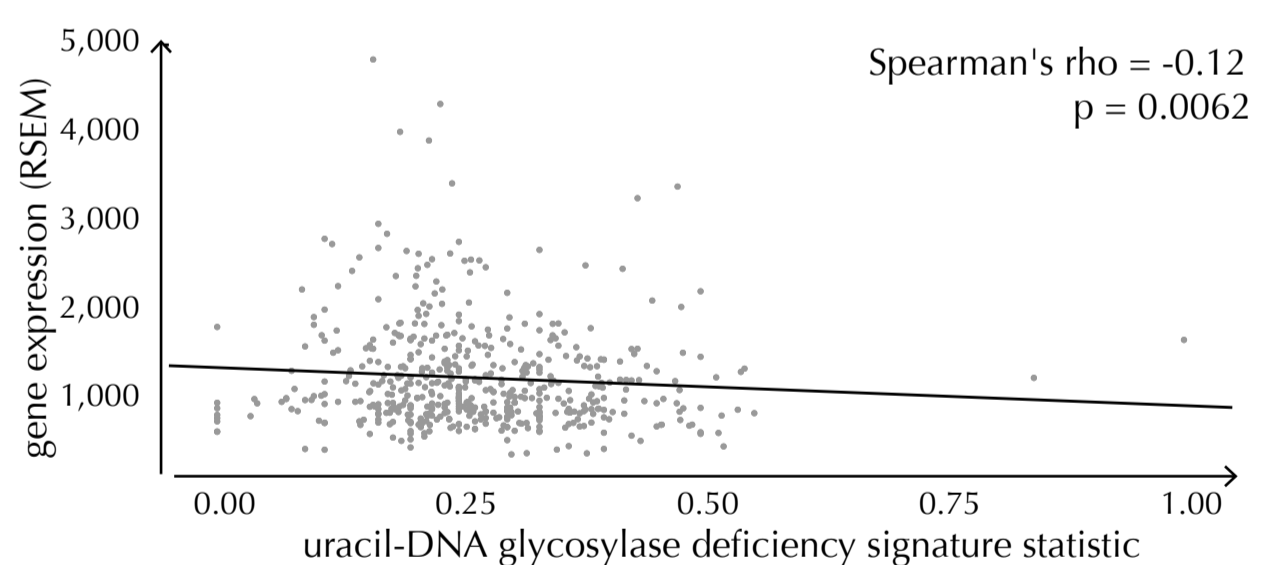
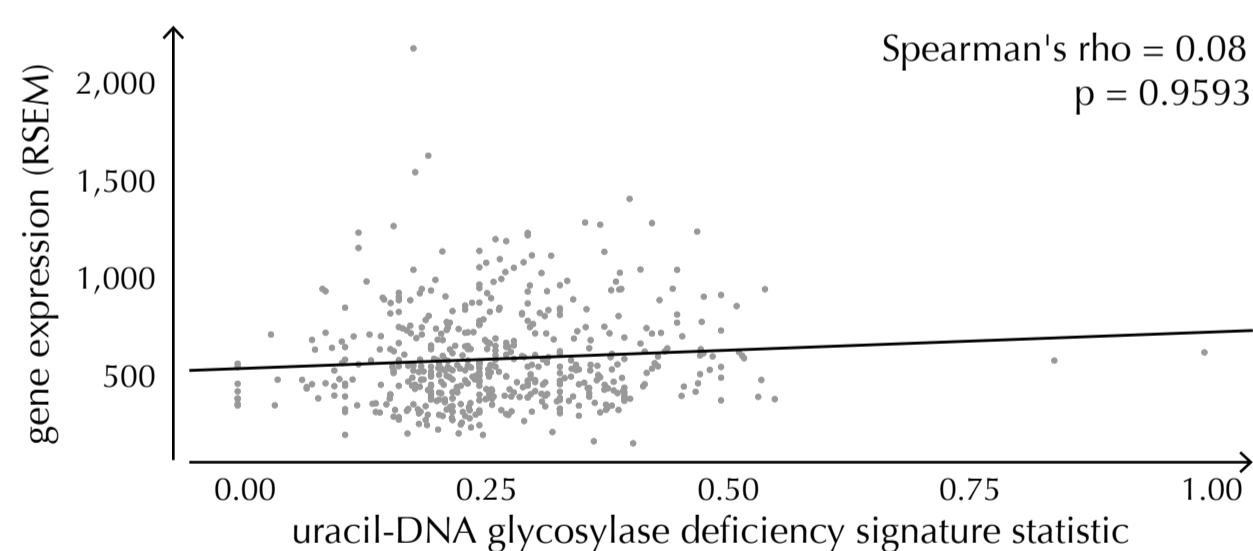
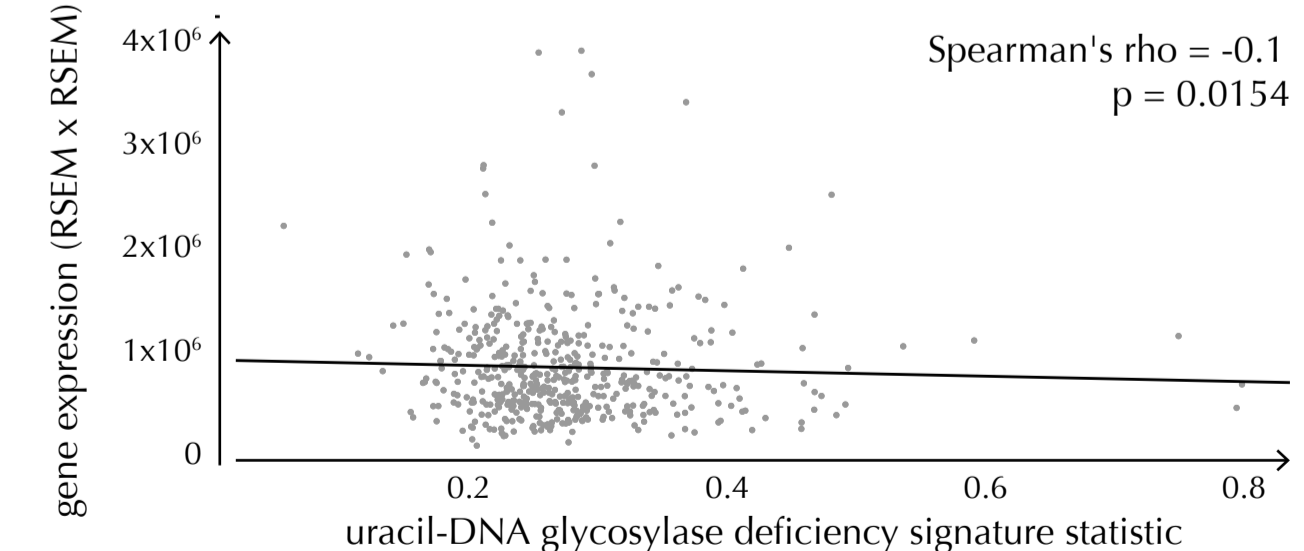
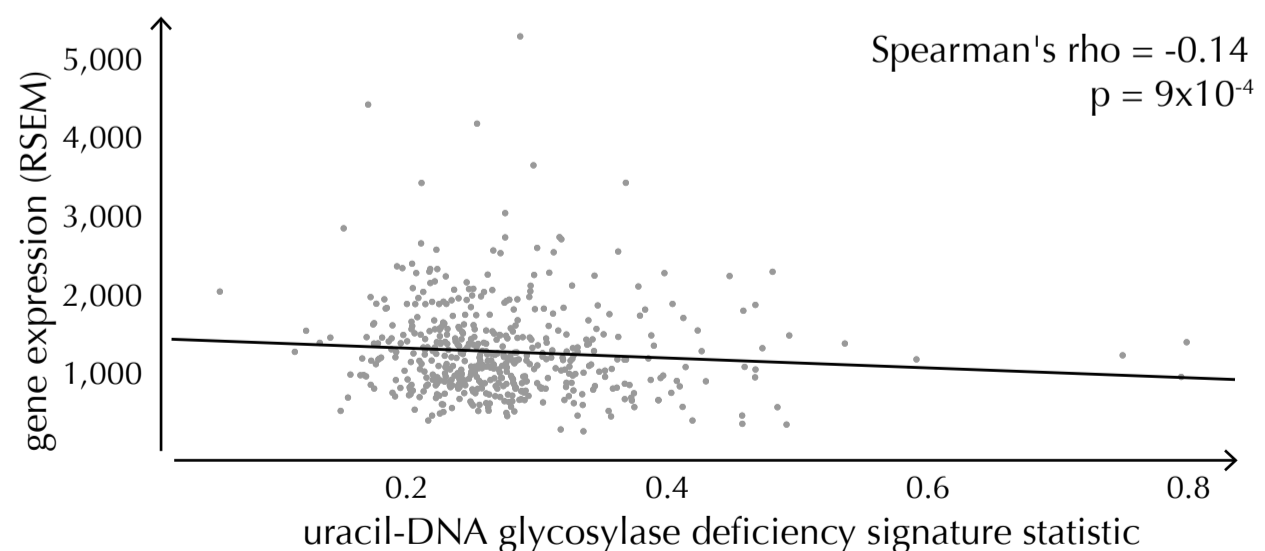
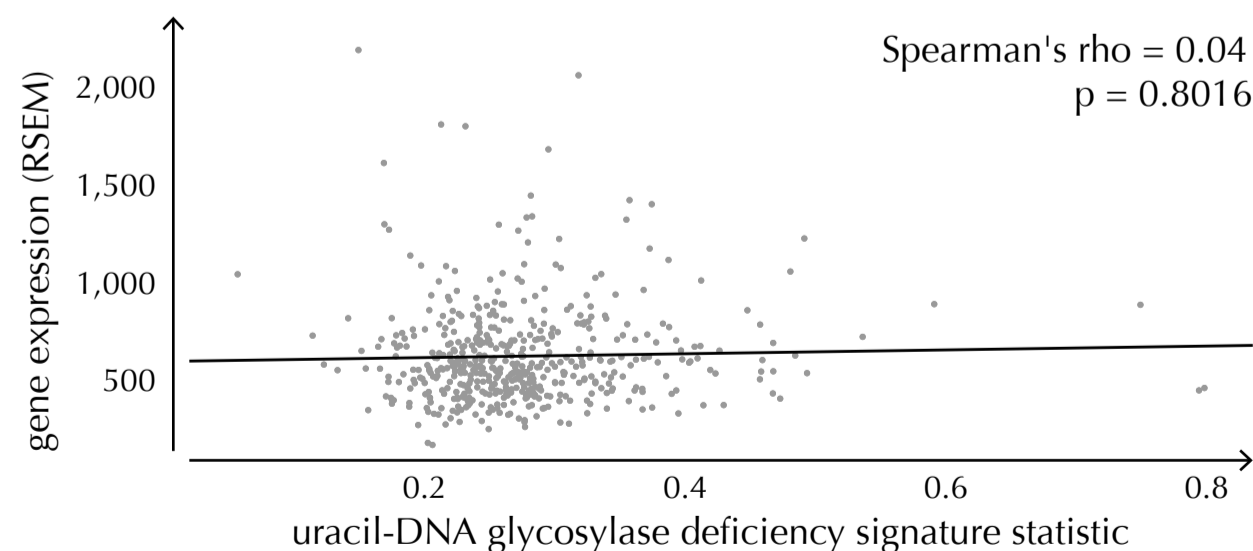
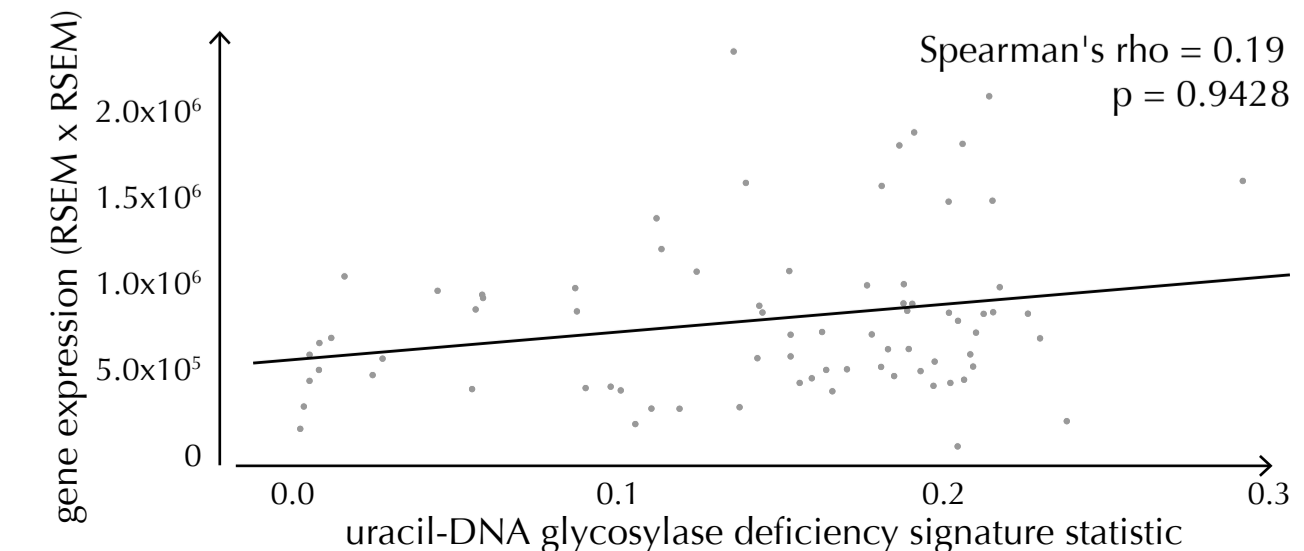
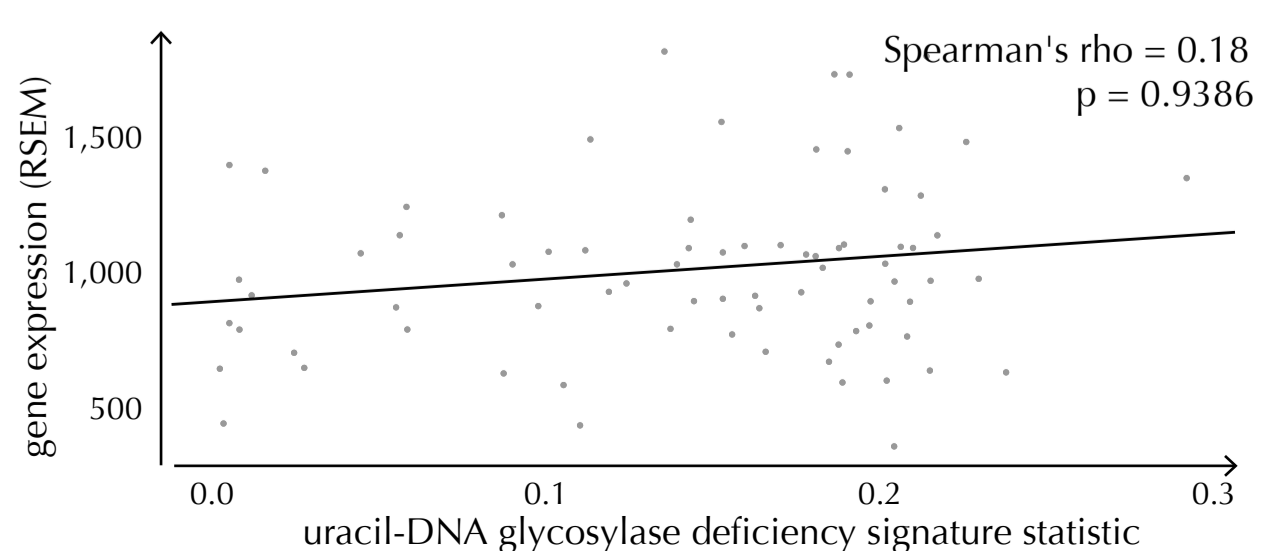
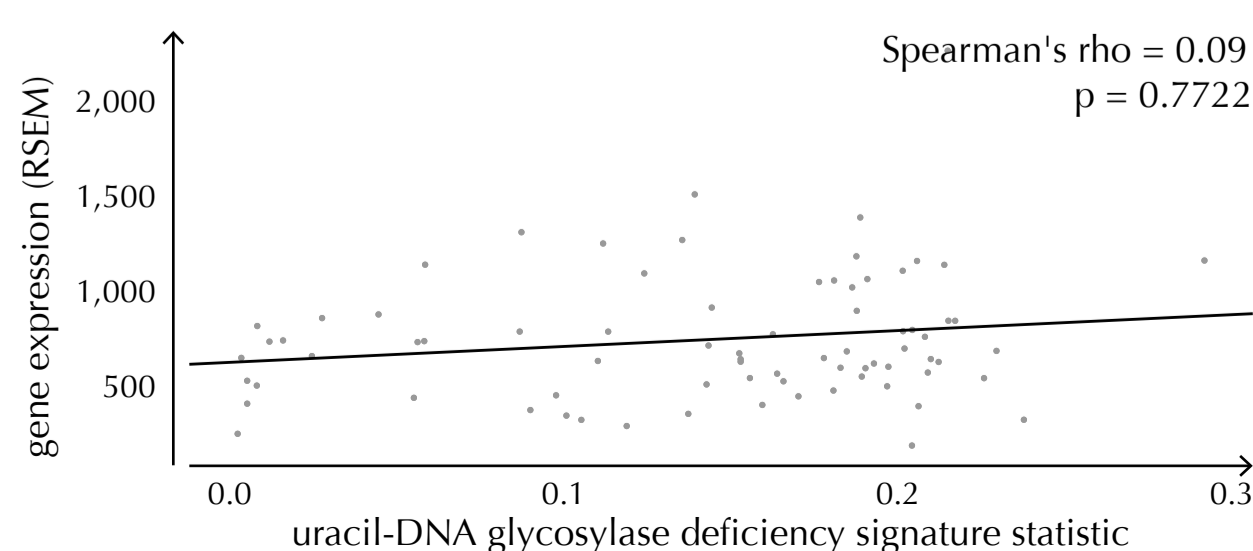
C 5' UTRs ($n = 18,220$)



rank	conservation		log likelihoods		number of mutations	
	gene name	q-value	gene name	q-value	gene name	q-value
1	WDR74	1.84×10^{-12}	WDR74	3.79×10^{-16}	WDR74	4.07×10^{-12}
2	C16orf59	3.66×10^{-10}	MRPL36	6.74×10^{-13}	MRPL36	2.02×10^{-9}
3	MRPL36	2.76×10^{-9}	MTG2	1.18×10^{-7}	C16orf59	2.12×10^{-8}
4	MTRNR2L13	3.89×10^{-5}	HLA-F	1.18×10^{-7}	MTRNR2L13	4.86×10^{-5}
5	HLA-F	1.02×10^{-4}	C16orf59	7.07×10^{-7}	MTG2	1.24×10^{-4}
6	MTG2	1.03×10^{-4}	ZNF717	9.47×10^{-7}	HLA-F	1.61×10^{-4}
7	NDUFB9	4.75×10^{-4}	NDUFB9	3.08×10^{-4}	EVI2A	5.34×10^{-4}
8	PRKAG1	4.75×10^{-4}	TBC1D12	3.08×10^{-4}	DHX16	1.46×10^{-3}
9	EVI2A	4.75×10^{-4}	MTRNR2L13	3.10×10^{-4}	ZNF717	1.46×10^{-3}
10	TBC1D12	5.74×10^{-4}	EVI2A	3.10×10^{-4}	PRKAG1	1.46×10^{-3}

gene name	WDR74	C16orf59	MRPL36	MTRNR2L13	HLA-F	MTG2	NDUFB9	PRKAG1	EVI2A	TBC1D12	DHX16	ZNF717
region size (bp)	839	173	166	941	384	26	84	487	548	110	323	600
observed mutations	30	12	18	20	19	5	6	11	9	6	10	27
expected mutations	2.4	0.6	0.9	3.7	2.2	0.1	0.3	1.4	0.7	0.3	1.1	4.4



SMUG1**UNG****SMUG1 x UNG****BLCA****LUAD****UCEC****ACC****BRCA****HNSC****SKCM**

