

1 **Ultra-deep sequencing differentiates patterns of skin clonal mutations associated with**
2 **sun-exposure status and skin cancer burden**

3
4 **Classification:** Biological Sciences, Genetics; Cancer risk; Early detection;

5
6 Lei Wei^{1,*§}, Sean R. Christensen^{2,*}, Megan Fitzgerald³, James Graham¹, Nicholas Hutson¹, Chi
7 Zhang⁴, Ziyun Huang⁵, Qiang Hu¹, Fenglin Zhan^{1,6}, Jun Xie⁷, Jianmin Zhang⁸, Song Liu¹, Eva
8 Remenyik⁹, Emese Gellen⁹, Oscar R. Colegio^{10,11}, Michael Bax¹⁰, Jinhui Xu¹², Haifan Lin¹³, Wendy
9 J. Huss^{14,*}, Barbara A. Foster^{14,*}, Gyorgy Paragh^{3,9,*§}

10 **Author affiliations:**

11 ¹Department of Biostatistics and Bioinformatics, Roswell Park Comprehensive Cancer Center,
12 Buffalo, NY

13 ²Department of Dermatology, Yale University School of Medicine, New Haven, CT

14 ³Department of Cell Stress Biology, Roswell Park Comprehensive Cancer Center, Buffalo, NY

15 ⁴School of Biological Sciences Center for Plant Science and Innovation, University of Nebraska,
16 Lincoln, NE

17 ⁵Department of Computer Science and Software Engineering, Penn State Erie, The Behrend
18 College

19 ⁶PET/CT center, The First Affiliated Hospital of USTC, Division of Life Sciences and Medicine,
20 University of Science and Technology of China, Hefei, Anhui, 230001, P.R. China

21 ⁷Department of Statistics, Purdue University, West Lafayette, IN

22 ⁸Department of Cancer Genetics and Genomics, Roswell Park Comprehensive Cancer Center,
23 Buffalo, NY

24 ⁹Department of Dermatology, Faculty of Medicine, University of Debrecen, Debrecen, Hungary

25 ¹⁰Department of Dermatology, Roswell Park Comprehensive Cancer Center, Buffalo, NY

26 ¹¹Department of Immunology, Roswell Park Comprehensive Cancer Center, Buffalo, NY

27 ¹²Department of Computer Science and Engineering, State University of New York at Buffalo

28 ¹³Yale Stem Cell Center, Yale University School of Medicine, New Haven, CT

29 ¹⁴Department of Pharmacology and Therapeutics, Roswell Park Comprehensive Cancer Center,
30 Buffalo, NY

31

32 * These authors contributed equally

33 § Corresponding authors: Lei.Wei@RoswellPark.org and Gyorgy.Paragh@RoswellPark.org

34

35 **Email addresses:**

36 LW: Lei.Wei@RoswellPark.org

37 SC: Sean.Christensen@Yale.edu

38 MF: Megan.Fitzgerald@RoswellPark.org

39 JG: james.graham@stonybrookmedicine.edu

40 NH: ndhutso@gmail.com

41 CZ: czhang5@unl.edu

42 ZH: zxh201@psu.edu

43 QH: qiang.hu@roswellpark.org

44 FZ: zhan209@hotmail.com

45 JX: junxie@purdue.edu

46 JZ: jianmin.zhang@roswellpark.org

47 SL: song.liu@roswellpark.org

48 ER: remenyik@med.unideb.hu
49 EG: emesegellen@med.unideb.hu
50 OC: Oscar.Colegio@RoswellPark.org
51 MB: Michael.Bax@RoswellPark.org
52 JX: jinhui@buffalo.edu
53 HL: haifan.lin@yale.edu
54 WJH: Wendy.Huss@RoswellPark.org
55 BAF: Barbara.Foster@RoswellPark.org
56 GP: Gyorgy.Paragh@RoswellPark.org

57

58 **Keywords:** ultraviolet light, clonal mutation, photocarcinogenesis, sun exposure, ultra-deep
59 sequencing, skin cancer risk

60

61 **Abstract**

62 Non-melanoma skin cancer is the most common human malignancy and is primarily caused by
63 exposure to ultraviolet (UV) radiation. The earliest detectable precursor of UV-mediated skin
64 cancer is the growth of cell groups harboring clonal mutation (CM) in clinically normal appearing
65 skin. Systematic evaluation of CMs is crucial to understand early photo-carcinogenesis. Previous
66 studies confirmed the presence of CMs in sun-exposed skin. However, the relationship between
67 UV-exposure and the accumulation of CMs, and the correlation of CMs with skin cancer risk
68 remain poorly understood. To elucidate the exact molecular and clinical effects of long-term UV-
69 exposure on skin, we performed targeted ultra-deep sequencing in 450 individual-matched sun-
70 exposed (SE) and non-sun-exposed (NE) epidermal punch biopsies obtained from clinically
71 normal skin from 13 donors. A total of 638 CMs were identified, including 298 UV-signature
72 mutations (USMs). The numbers of USMs per sample were three times higher in the SE samples
73 and were associated with significantly higher variant allele frequencies (VAFs), compared with
74 the NE samples. We identified genomic regions in *TP53*, *NOTCH1* and *GRM3* where mutation
75 burden was significantly associated with UV-exposure. Six mutations were almost exclusively
76 present in SE epidermis and accounted for 42% of the overall difference between SE and NE
77 mutation burden. We defined Cumulative Relative Clonal Area (CRCA), a single metric of UV-
78 damage calculated by the overall relative percentage of the sampled skin area affected by CMs.
79 The CRCA was dramatically elevated by a median of 11.2 fold in SE compared to NE samples.
80 In an extended cohort of SE normal skin samples from patients with a high- or low- burden of
81 cutaneous squamous cell carcinoma (cSCC), the SE samples in high-cSCC patients contained
82 significantly more USMs than SE samples in low-cSCC patients, with the difference mostly
83 conferred by mutations from low-frequency clones (defined by $VAF \leq 1\%$) but not expanded clones
84 ($VAF > 1\%$). Our studies of differential mutational features in normal skin between paired SE/NE
85 body sites and high/low-cSCC patients provide novel insights into the carcinogenic effect of UV

86 exposure, and indicate that CMs might be used to develop novel biomarkers for predicting cancer
87 risk.

88 **Significance statement:**

89 In UV radiation exposed skin, mutations fuel clonal cell growth. We established a sequencing-
90 based method to objectively assess the mutational differences between sun-exposed (SE) and
91 non-sun-exposed (NE) areas of normal human skin. Striking differences, in both the numbers of
92 mutations and variant allele frequencies, were found between SE and NE areas. Furthermore, we
93 identified specific genomic regions where mutation burden is significantly associated with UV-
94 exposure status. These findings revealed previously unknown mutational patterns associated with
95 UV-exposure, providing important insights into UV radiation's early carcinogenic effects.
96 Additionally, in an extended cohort, we identified preliminary association between normal skin
97 mutation burden and cancer risk. These findings pave the road for future development of
98 quantitative measurement of subclinical UV damage and skin cancer risk.

99 **Background**

100 Ultraviolet (UV) light is responsible for over 5 million cases of skin cancer annually in the US,
101 which is more human malignancies than all other environmental carcinogens combined^{1,2}. In
102 mammals, nucleotide excision repair eliminates UV-mediated DNA lesions, but this mechanism
103 of repair is error prone resulting in frequent mutations³. The preferential location of UVB induced
104 DNA lesions results in a specific pattern of so-called UV signature mutations at dipyridine sites
105 (C>T, CC>TT)⁴. In most skin cancers, including cutaneous squamous cell carcinoma (cSCC), the
106 burden of UV signature driver mutations is high^{4,5}. While some cSCC arise from visible
107 precancerous lesions known as actinic keratoses (AKs), many cSCC arise in apparently "normal"
108 skin areas from precursors that are clinically invisible⁶. Therefore, clinically visible precursors are
109 an ominous sign but not a sensitive early measure of photocarcinogenesis.

110 *TP53* mutations are among the most common driver mutations in cSCC, and are also detected
111 by immunohistochemistry in aged normal skin^{7,8}. These UV-induced *TP53* mutations facilitate
112 clonal expansion of cells harboring them and therefore behave as early clonal mutations (CMs)⁹.
113 For two decades *TP53* mutant keratinocyte cell clones were considered the earliest
114 manifestations of skin carcinogenesis^{7,8,10}. Because p53 clonal immunopositivity could not be
115 efficiently quantified in human skin, detection of mutant *TP53* for assessment of
116 photocarcinogenesis in clinical dermatology practice has been unattainable. The low relative
117 abundance of clonal DNA previously limited efficient detection of early mutated cell groups.
118 However, with improved high throughput sequencing technology we have finally reached the
119 lower end of this threshold and efficient detection of rare mutations in normal tissue is becoming
120 feasible in recent studies by others and us using deep bulk sequencing or single cell DNA
121 sequencing¹¹⁻¹⁶. In exploratory analyses, CMs were found to be abundant in clinically normal skin
122 from sun-exposed sites in *NOTCH1*, *NOTCH2*, *FAT1* and several other genes besides *TP53*¹².
123 Prior attempts to establish a quantitative method for assessing photodamage and skin cancer risk
124 had limited success^{17,18}. A method that enables quantitative evaluation of early photodamage is
125 expected to help optimize personalized sun-protective measures and may also serve as a tool for
126 assessing the need and efficacy of early preventative treatment interventions.

127 In the current work we developed an ultra-deep sequencing-based method to identify CMs in
128 clinically normal epidermis and show differences in CMs between sun-exposed and non-sun-
129 exposed skin areas. We then correlated CMs with skin cancer burden in another independent
130 cohort of cSCC patients and found mutational features in normal skin are significantly associated
131 with cancer risk burden.

132

133 **Methods**

134 **Samples:**

135 A total of 464 normal human skin samples were collected from 13 Caucasian post-mortem donors
136 over the age of 55 years using Roswell Park's Rapid Tissue Acquisition Program under a Roswell
137 Park approved IRB protocol within 24 hours of death from frequently sun-exposed (SE) sites (left
138 dorsal forearm) and non-sun-exposed (NE) sites (left medial buttock). Exclusion criteria included
139 any visible skin abnormalities in the tissue areas. Eligible donors were identified and clinically
140 normal appearing skin was harvested. Skin samples were kept in tissue preservation medium,
141 Belzer UW cold storage solution (Bridge to Life, USA) at 4°C until processed. All samples that
142 could be processed within 36 hours or less after death were included in the study. The mean age
143 of the donors was 72.3 years (SD: ± 8.2 years; range 60-80 years). The male to female gender
144 ratio was 7:6, and 12/13 donors had no history of skin cancer.

145 The adipose tissue was removed from each human skin sample using sterile scissors. The
146 samples were cut into strips wide enough to harvest 6 mm punches. The epidermis was separated
147 from the dermis by placing the strips in tubes containing 10 ml of 5U/ml Dispase II (Stem Cell
148 Technologies, USA) and incubated at 4°C overnight and at 37°C for 2-3 hours. After Dispase
149 digestion the specimens were placed in a petri dish containing a small amount of 1x DPBS
150 (Corning, USA) and using sterile tweezers, the epidermis was carefully removed from the dermis.
151 Using disposable biopsy punches, 1, 2, 3, 4 and 6 mm diameter epidermal pieces were taken
152 from the epidermal sheets and punched epidermal pieces were placed into a sterile 1.5 mL vials.
153 In addition to the epidermal punches, large bulk pieces of dermis were also removed from the
154 skin samples using a disposable #15 blade and placed into a sterile 1.5 mL vial for use as a
155 germline control.

156 For the extended cohort of the study, 20 human skin samples were obtained in a de-identified
157 manner from 8 undergoing surgery for cSCC. The mean age of the donors was 77.9 years (SD:
158 ± 12.3 years; range 54-92 years). The male to female gender ratio was 1:1. The study was granted
159 exemption by the Yale University Human Investigation Committee (Protocol 1509016421). All
160 individuals had biopsy-confirmed cSCC that was completely excised by Mohs micrographic
161 surgery with intraoperative histologic verification of clear surgical margins. Immediately following
162 excision of cSCC, adjacent normal skin was excised to facilitate surgical repair and samples for
163 sequencing were immediately harvested. From each individual, two skin samples at a fixed linear
164 distance from the cSCC were obtained from the adjacent, sun-exposed, normal skin. One sample
165 was obtained at a distance of 1mm from the cSCC surgical margin, and one at a distance of 6mm
166 from the surgical margin. From four patients, a tumor sample from grossly visible cSCC was also
167 obtained at the time of surgery. All samples were obtained with a 2mm punch biopsy to a depth
168 of approximately 1mm, including epidermis and superficial dermis.

169 **DNA isolation:**

170 DNA samples from the primary cohort were extracted using Purelink™ Genomic DNA mini kit
171 (Invitrogen, USA). Epidermal samples were digested using Proteinase K at 55°C heating block
172 overnight following the manufacturers recommendations. For the extended cohort of samples,
173 skin biopsies were similarly digested using Proteinase K and DNA was purified with phenol-
174 chloroform extraction and ethanol precipitation. DNA was eluted with 28 μ L of Molecular Biology
175 Grade Water (Corning, USA) for 1 and 2 mm punches or 36 μ L of Molecular Biology Grade Water
176 for 3, 4, and 6 mm punches. The isolated genomic DNA was stored at -20°C and the DNA
177 concentration of each extraction was measured using a Qubit fluorometer or Quanti-iT PicoGreen
178 kit (Invitrogen, USA).

179 **Ultra-deep Targeted Sequencing:**

180 The sequencing libraries were generated using the TruSeq Custom Amplicon kit (Illumina, USA)
181 using 10-50 ng of gDNA. Amplicons of ~150bp (primary cohort) or ~250bp (extended cohort) in
182 length were designed using Illumina Design Studio Software. Custom oligo capture probes that
183 flank the regions of interest were hybridized to the gDNA. A combined extension/ligation reaction
184 completed the region of interest between these flanking custom oligo probes. PCR was then
185 performed to add indices and sequencing adapters. The amplified final libraries were cleaned up
186 using AmpureXP beads (Beckman Coulter). Purified libraries were run on a TapeStation
187 DNA1000 screentape chip to verify desired size distribution, quantified by KAPA qPCR (KAPA
188 Biosystems) and pooled equal molar in a final concentration of 2 nM. Pooled libraries were loaded
189 on an Illumina HiSeq Rapid Mode V2 flow cell following standard protocols for 2x100 cycle
190 sequencing (primary cohort), or Illumina NextSeq for 2x150 cycle sequencing (extended cohort).

191 **Bioinformatics analysis:**

192 High quality paired-end reads passing Illumina RTA filter were initially processed against the NCBI
193 human reference genome (GRCh37) using public available bioinformatics tools ^{19,20}, and Picard
194 (<http://picard.sourceforge.net/>). The coverage quality control required at least 80% of the targeted
195 region covered by a minimum of 1,000X coverage. Putative mutations, including single nucleotide
196 variants (SNVs) and small insertions/deletions (Indels), were initially identified by running variation
197 detection module of Strelka²¹ on each SE or NE epidermis sample paired with the matched dermal
198 sample. From the detected SNVs, dinucleotide variants (DNV) or cluster of single nucleotide
199 variants (CSNV) were recognized by running Multi-Nucleotide Variant Annotation Corrector (MAC)
200 ²² on the original sequences. The putative mutations detected from all samples were consolidated
201 into a list of unique mutations. Every unique mutation was re-visited in all samples to calculate
202 the numbers of mutant/wildtype reads, as well as variant allele frequency (VAF) in each sample
203 as previously described ¹³.

204 To distinguish mutations from background errors, we modelled each mutation's background
205 error rate distributions using VAFs from all control (dermal) samples. For each mutation, we
206 started by fitting a *Weibull* distribution to VAFs from all control samples following a previously
207 published method²³, then every SE or NE epidermal sample's VAF was compared to the fitted
208 distribution. A positive sample was defined as the sample's VAF of a mutation was significantly
209 above background ($p < 0.05$, after Bonferroni correction). In the extended cohort where the control
210 samples were not available, we adapted a dynamic control strategy, based on the assumption
211 that any somatic mutation cannot be recurrent in more than 10% of all samples at the same site.
212 In the previous primary cohort, all recurrent mutations were within 5% of all samples. For each
213 potential mutation, we first cluster the VAFs of the mutation in all samples. Subsequently started
214 from the cluster with lowest VAF, we transferred all samples of each cluster to the control cohort
215 until at least 90% of all samples are in the control cohort. After mutation calling, all identified
216 mutations including SNVs, DNVs, CSNVs and Indels were annotated using a customized program
217 with NCBI RefSeq database.

218 Cumulative Relative Clonal Area (CRCA), defined as the overall percentage of biopsied skin
219 area covered by UV-signature mutations (USMs) in a patient skin punch, was calculated as
220 following:

$$221 \quad CRCA = \frac{\sum_{i=1}^n \sum_{j=1}^{m_i} (\pi r_i^2 * 2VAF_j)}{\sum_{i=1}^n \pi r_i^2}$$

222 with n = the total number of punches collected in the patient; r_i = the size (radius) of each punch;
223 m_i = the number of mutations in punch i ; VAF_j = the variant allele fraction of a specific mutation j .
224 Here, the calculation of CRCA was based on the assumption that all mutations occur in one
225 chromosome of regular diploid genomic regions. Additionally, although we did not consider the

226 situation when multiple mutations occur in the same cell, we did identify mutations that occur on
227 the same reads and combined them into one mutation using MAC²².

228 **Statistics:**

229 The overall mutation numbers and VAFs between two groups, including SE and NE in the primary
230 cohort, and the high- and low- cSCC burden in the extended cohort, were evaluated using a
231 Wilcoxon test. Group-specific markers, including mutations, genes, regions and signatures were
232 identified using a Fisher's exact test where the two variables in the contingency table were the
233 samples' sun-exposure status (SE vs NE, in cohort #1) or cSCC burden (high vs low, in cohort
234 #2) and mutational status. Multiple testing correction was implemented using the FDR approach
235 as indicated.

236 **Results**

237 **Ultra-deep sequencing of epidermal samples using customized focused panels**

238 To generate a focused sequencing panel targeting the most commonly mutated sequences in
239 normal human skin, we selected an area of focus based on a previous dataset¹². All previous
240 mutations were assigned to 100-bp genomic segments. After sorting the segments by number of
241 mutations, we designed a panel to capture the top 55 most frequently mutated segments from 12
242 genes (5.5 kb in total, **Table S1**). The majority (65%) of the targeted segments came from the
243 following 3 genes: *NOTCH1*, *NOTCH2*, and *TP53*. When summarized by coding regions, 79% of
244 the targeted segments lie in protein-coding regions, and the remaining segments were mostly in
245 introns. In the previous dataset¹², 87% of the samples harbored at least one mutation within this
246 panel. Thus, as designed, this panel captured the most frequently mutated genomic regions in
247 sun-exposed skin, and was highly focused for efficient deep-sequencing to identify low-frequency
248 mutations.

249 The primary cohort was sequenced using the focused panel in two batches. We first
250 sequenced a discovery cohort of 374 human skin samples from 13 post-mortem donors: 360
251 epidermal samples, equally acquired from both sun-exposed (SE) and non-sun-exposed regions
252 (NE) using 1 mm, 2 mm, 3 mm, 4 mm or 6 mm punch sizes. From the same 13 donors, DNA
253 from bulk NE dermis (n=14, 1 donor contributed 2 samples) was isolated for germline controls.
254 After initial analysis to determine the optimal punch size, we then tested a separate validation
255 cohort of 90 epidermal samples from 9 of the 13 donors using the most effective punch size (2
256 mm, as detailed in results “Optimization of punch size for USM detection”). In total, the dataset
257 contains 464 samples: 225 SE, 225 NE, and 14 dermal samples as controls (**Table 1**) from 13
258 individuals. After sequencing, 85% of samples reached a minimum of 10,000X coverage in at
259 least 80% of the targeted region. The median of average coverage across all samples was
260 64,730X (**Table S2a**), with only one sample exclusion (NE sample) due to sequencing failure.
261 This unique design of ultra-deep sequencing from individual matched SE/NE samples enabled us
262 to discriminate between the mutational profiles of SE and NE skin samples.

263 To better define the clinical relevance of CMs, we sequenced an extended cohort of sun-
264 exposed skin samples from human patients with cSCC. Twenty 2mm punch biopsy specimens
265 were obtained from surgically excised skin from 8 individuals, including 16 normal skin samples
266 and 4 samples of cSCC. For this extended cohort, a custom sequencing panel was designed to
267 encompass the complete protein-coding region of 12 genes with frequently reported mutations in
268 UV-exposed skin (*NOTCH1*, *NOTCH2*, *NOTCH3*, *TP53*, *CDKN2A*, *BRAF*, *HRAS*, *KRAS*, *NRAS*,
269 *KNSTRN*, *FAT1*, and *FGFR3*), and 1 control gene without expected functional significance in skin
270 (*VHL*). This sequencing panel encompassed 59.5 kb. After sequencing, all samples have at least
271 80% of the targeted region covered by a minimum of 10,000X coverage. The median value of
272 average coverages across all samples was 47,158X (**Table S2b**). This extended cohort from
273 cSCC patients would allow us to correlate the features of CMs to patient clinical outcomes.

274 **Delineate the mutational patterns associated with UV exposure**

275 To identify the mutations solely caused by UV exposure, we characterized the mutational profiles
276 of individual-matched SE/NE epidermal samples. Additionally, we compared the epidermal
277 samples to patient-matched dermal samples followed by an in silico error suppression to remove
278 germline polymorphisms and low-frequency technical artifacts. Dinucleotide and other complex
279 mutations were identified by re-visiting the raw reads using a program that we previously
280 developed²². Altogether, a total of 638 mutations were identified, predominantly single nucleotide
281 variants (SNVs, n = 614 or 96.2%) or dinucleotide variants (DNVs, n = 20 or 3.1%) (**Table S3**).
282 The median variant allele frequency (VAF) of all mutations was 2.1% (range 0.1% - 36.6%), and
283 only 3% mutations reached a VAF greater than 10%.

284 Among the 55 targeted genomic segments, mutations were detected in 50 segments with an
285 average of 7.1 and 4.7 mutations per segment in SE and NE samples, respectively (**Figure 1a**).
286 Two segments were significantly (FDR $p < 0.001$) associated with UV-exposure status,
287 approximately corresponding to *TP53* amino acids 227-261 ("*TP53-3*", mutations in SE/NE = 38/0)
288 and *NOTCH1* p.449-481 ("*NOTCH1-9*", mutations in SE/NE = 30/4). Interestingly, mutations in
289 an adjacent region in *NOTCH1* p.419-449 ("*NOTCH1-10*") were not associated with UV exposure
290 (mutations in SE/NE=48/40), even though "*NOTCH1-10*" was the most frequently mutated
291 segment in the current study. Additionally, mutations were marginally enriched in SE samples
292 (FDR $p < 0.1$) in three other segments: two in *NOTCH1* ("*NOTCH1-14*" and "*NOTCH1-19*") and
293 one in *GRM3* ("*GRM3-2*"). On the gene level, mutations in SE samples were only significantly
294 enriched in *TP53* (FDR $p < 0.001$), and marginally significant in *GRM3* (FDR $p < 0.1$). Overall, the
295 numbers of mutations in SE samples were 6.3 times higher than NE samples in *TP53*, and 4.3
296 times in *GRM3* (**Figure 1b**). Mutations identified in nine other genes did not exhibit significant
297 association with sun-exposure status either on the gene- or segment- level: *NOTCH2*, *ARID1A*,

298 *SALL1*, *SCN1A*, *ERBB4*, *FAT4*, *FGFR3*, *ADGRB3* and *PPP1R3A*. These findings indicate a
299 highly genomic-region-specific pattern of the accumulation of UV-induced somatic mutations.

300 We next investigated potential hotspots and mutations associated with UV-exposure. After
301 sorting all mutations by their genomic locations, one specific region in *TP53* (p.217-280),
302 appeared to be “mutation exempt” in comparison to surrounding regions in NE samples. In
303 contrast, this region was highly mutated in SE samples (**Figure 2a**). We reanalyzed a recent study
304 involving RNASeq of both SE and NE normal skin samples¹¹, and found four mutations in this
305 region, all from SE samples (**Table S4**). To identify mutations associated with UV exposure, we
306 focused on highly recurrent mutations (present in ≥ 5 samples, $n = 18$). By comparing the
307 frequency in SE and NE skin samples, we identified six mutations significantly enriched in SE
308 samples: *TP53* R248W, *NOTCH1* P460L, *NOTCH1* S385F, *NOTCH1* E424K, *TP53* G245D and
309 *NOTCH1* P460S, and nearly all of them were exclusively found in SE samples (FDR $p < 0.05$,
310 **Figure 2b**). No mutation was significantly enriched in NE samples. Five of the six SE-enriched
311 mutations were found in both discovery and validation cohorts, indicating they were unlikely to be
312 caused by batch-effect. Unexpectedly, one specific mutation (*NOTCH1* E424K) was associated
313 with significantly elevated VAFs (median = 10%, $p < 0.001$, Wilcoxon test), about five-fold higher
314 than other mutations (median VAF = 2.1%, **Figure 2a, 2b**). Through protein structure modelling
315 (**Figure 2c**), we found that the *NOTCH1* E424K mutation is predicted to disrupt the binding of
316 *NOTCH1* to delta-like canonical ligand 4 (*DLL4*), a negative regulator of the Notch signaling
317 pathway¹¹. By prohibiting formation of a salt bridge between *NOTCH1* E424 and *DLL4* K189/R191,
318 the mutation E424K creates a repulsive force that inhibits *DLL4* binding²⁴. Based on the biological
319 role of *DLL4* and *NOTCH1*, the *NOTCH1* E424 mutation is expected promote epithelial
320 proliferation^{25,26}. The overall prevalence of the *NOTCH1* E424K mutation in our dataset is 2.7%.
321 For comparison, in GENIE cBioPortal²⁷, *NOTCH1* E424K is mutated in 1.3% of cutaneous SCCs,
322 0.04% in melanomas, and is rarer in other cutaneous or non-cutaneous malignancies (**Table S5**).

323 **UV-signature mutations exclusively account for the elevated mutation burdens in SE skin**

324 We next intercorrelated the identified mutations with previously known UV-signature mutations
325 (USMs), i.e., C>T transition at dipyrimidines⁴. Among all 638 mutations in SE and NE samples,
326 298 were USMs. Of these 298 USMs, 76% were present in SE samples. USMs were significantly
327 enriched in SE compared to NE samples (n = 226 and 72, respectively, p<0.001, Fisher's exact
328 test). Especially among the high-VAF mutations, 18 of 19 mutations with VAFs above 0.1 were
329 from SE samples, and most (13 of 18) were USMs. Conversely, non-UV-signature mutations
330 (NUSMs) were present approximately equally (n= 159 and 181, ns, Fisher's exact test) in SE and
331 NE skin types (**Figure 3a**), indicating that these mutations may not be directly associated with
332 UV-exposure.

333 To explore specific community enrichment patterns in different mutational function groups, we
334 classified all 638 mutations into four effect-groups: nonsense, missense, silent and noncoding.
335 Inside each effect-group, we correlated the mutational properties (USM vs NUSM) with the
336 matched samples' sun-exposure statuses (SE vs NE) (**Figure 3b**). Significant enrichment of
337 USMs were observed in two of four effect-groups by Fisher's exact test: nonsense (FDR p<0.05)
338 and missense (FDR p<0.001). Specifically, nonsense mutations were 9 times more frequently
339 occurring in SE skins than in NE skins, and similarly enriched by 4.2 times for missense mutations.
340 These findings indicate that the mutations initiated by UV radiation are further selected by the
341 host system or inter-clonal competition²⁸, in which the mutations with functional impacts give the
342 host clone greater fitness.

343 **Quantification of UV-induced DNA damage level by UV-signature mutations**

344 We next investigated the feasibility of using CMs to quantify UV-induced DNA damage. This was
345 based on the hypothesis that SE samples harbor more CMs and are associated with higher VAFs
346 compared to NE samples. Since our analyses indicated NUSMs were not correlated with UV

347 exposure, only USMs were used for quantifying UV-induced DNA damage. To avoid the potential
348 bias introduced by different punch sizes, initially only the most abundant size of 2 mm (n = 90 and
349 89, SE and NE, respectively) (**Figure 3c**) was analyzed. A three-fold difference was observed in
350 the average USMs per sample between SE (mean = 1.2) and NE (mean = 0.4), which was
351 significantly higher ($p < 0.001$, Wilcoxon test). Multiple USMs were found in 33% of SE samples
352 but only 9% of NE samples (**Table S6**). Additionally, the identified USMs had significantly higher
353 VAFs in SE (mean = 3.7%) than NE (mean = 2.1%) samples, ($p < 0.001$, Wilcoxon test), indicating
354 the presence of larger clones in SE samples (**Figure 3d**). We further extended the analysis to
355 include all punch sizes, and found the pattern was consistent with 34% of SE and only 6% of NE
356 samples having multiple USMs and three-fold higher average USMs per sample in SE (1.0) than
357 NE (0.3) samples ($p < 0.001$, Wilcoxon test). These findings of increased USMs and elevated
358 VAFs in SE than NE skin would then serve as the cornerstones for the quantification of UV-
359 induced DNA damage.

360 In order to overcome the heterogeneity between samples, we developed Cumulative Relative
361 Clonal Area (CRCA) as a single metric to assess the overall patient-level burden of CMs. The
362 CRCA was defined as the overall percentage of biopsied skin area covered by USMs in a patient
363 skin punch, which account for both the number of USMs and their VAFs (**Figure 3e**). It is worth
364 mentioning that our data did not allow us to distinguish whether mutations occurred independently
365 or were present in the same clone. Hence, CRCA does not provide an exact measure of the
366 mutated cell population, but rather serves as an index of the mutation burden in the sampled area.
367 To minimize the potential chance for repeated counting of co-occurring mutations in the same
368 cells, co-occurring mutations were identified, primarily dinucleotide CC>TT mutations, and
369 consolidated. When counted separately by sun-exposure status, the median CRCA across the 13
370 patients was 6.1% (range 1.4-14.2%) in SE and 1.4% (range 0.1-4.0%) in NE sites. On individual
371 patient level, the CRCAs were higher in SE than the matched NE skin in all patients, with the

372 average ratio of 11.2-fold higher (range = 1.4 - 55.0-fold). These CRCAs were calculated using
373 only USMs. If all CMs were included, the CRCA would be only 2.2-fold higher (range = 0.8 - 5.6-
374 fold) in SE than NE skin (data not shown). Based on these results, CRCA may have the potential
375 to be used as an objective measurement of the level of UV-induced DNA damage.

376 **The effect of punch size on USM detection**

377 In the discovery cohort, we sought to evaluate different punch sizes to determine the most efficient
378 one for detecting USMs. Theoretically, although larger punches likely contain more clones, they
379 tend to become less effective for detecting smaller clones due to a dilutional effect by other clones
380 harboring no or different mutations (**Figure 4a**). Overall across all five punch sizes, USMs were
381 detected in 54% of the SE, which was significantly higher than the 21% of the NE ($p < 0.001$,
382 Fisher's exact test). Between different punch sizes, 2 mm punches were found to have the highest
383 positive rate of 64%, and with the most significant difference between SE and NE samples ($p <$
384 0.0001 , **Figure 4b**). Thus, only 2mm punches were collected in the 90-sample validation cohort
385 and the extended cohort from cSCC patients. In the validation cohort, similarly, we found the SE
386 samples had higher numbers of USMs and the positive rate of USMs (69%) was similar to the
387 discovery cohort (64%).

388 When combining the discovery and validation cohorts, the SE samples had the highest positive
389 rate of 67% for USMs in 2 mm samples and were significantly higher than NE samples ($p < 0.001$),
390 followed by 60% in 4 mm ($p < 0.05$), and 54% in 3 mm ($p < 0.05$). Interestingly, the USM positive
391 rates were relatively lower in the largest punch size of 6 mm (53%) and the smallest 1 mm (36%).
392 In all NE samples, positive USM rates ranged from 17-30% (**Figure 4c**). Moreover, the punch
393 size also affected the detected VAFs of the mutations. Specifically, in SE samples, larger punches
394 were associated with smaller VAFs. The VAFs' standard deviation was the highest in 1 mm
395 punches (8.9%) and decreased with punch size: 2 mm (4.3%), 3 mm (2.8%), 4 mm (2.6%) and 6
396 mm (1.7%). This trend, between VAF range and punch size, was not present in NE samples

397 **(Figure 4d)**. These results suggested that the most effective punch size in detecting USMs under
398 the current sequencing condition was 2 mm.

399 **Mutation nucleotide contexts enriched with UV-exposure**

400 We next assessed the enrichment of different mutation nucleotide contexts in SE skin. The
401 mutation nucleotide contexts were defined by each SNV's trinucleotide and DNV's dinucleotide
402 contexts. A total of 83 contexts were identified from current mutations, including 13 contexts
403 matching to previously described USMs⁴. None of the remaining 70 non-USM contexts were
404 enriched in SE or NE samples **(Figure 5a)**. The 13 previously defined USM contexts were not
405 equally enriched in SE samples. After multiple test correction, only 5 of the 13 contexts were
406 significantly enriched in SE samples (FDR $p < 0.05$), including the dinucleotide CC>TT context,
407 which was exclusively found in SE samples **(Figure 5b)**. The most significant mutation context
408 enriched in SE samples was T[C>T]C (FDR $p = 0.00013$), which was in consonance with the
409 previously defined "Mutational Signature #7" in skin cancers²⁹. The remaining eight UV-signature
410 contexts were not significantly enriched in SE samples. Of particular note, G[C>T]C, which was
411 the most abundant context by total number of mutations, appeared to be equally presented in SE
412 and NE skin samples and therefore not associated with sun-exposure.

413 **Clonal mutations are correlated with cSCC burden**

414 To define the clinical significance of CMs and investigate the potential association with skin cancer
415 risk, we sequenced an extended cohort of 20 samples (16 SE normal skin and 4 cSCC) from eight
416 patients with cSCC using a 59.5 kb customized panel as described above. Four individuals
417 (including 8 normal skin samples and 2 cSCC samples from face, scalp, and arm) had a low
418 burden of skin cancer with only a single diagnosis of cSCC and few AKs (low-cSCC). Four
419 individuals (including 8 normal skin samples and 2 cSCC samples from face, hand, and lower leg)
420 had a high burden of skin cancer with severe UV damage, multiple prior cSCC (range 3-10) and

421 many AKs (high-cSCC). Low-cSCC and high-cSCC patients were matched for age (mean age
422 76.5 and 79.3, respectively). Normal skin samples were all sun-exposed, and were obtained a
423 linear distance of either 1mm or 6mm from the clear surgical margin of the excised cSCC, allowing
424 for analysis of CMs arising in skin subjected to carcinogenic UV radiation. Visible AKs were not
425 present in normal skin samples. A total of 535 somatic mutations were identified (**Table S7**), with
426 a median VAF of 1.2%. Only 15 mutations had VAF greater than 10%, most of which (10 of 15)
427 were from the cSCC tumor samples (**Figure 6a**). The median numbers of mutations per sample
428 in each group were 22 and 17.5 for the high- and low- cSCC normal skin samples (marginally
429 significant, $p=0.078$, Wilcoxon), and 41.5 for the cSCC samples. The overall mutation rates in
430 normal skin were 0.45 and 0.29 mutations per MB, in high- and low-cSCC patients, respectively.
431 The latter was comparable to the rate of SE normal skin of non-cancer patients in the primary
432 cohort (0.31 mutations per MB), despite the technical differences between the two cohorts such
433 as sequencing depth, targeted regions and punch sizes.

434 The frequently mutated genes in normal skin (more than two mutations per gene on average)
435 included *FAT1*, *NOTCH1*, *NOTCH2*, *NOTCH3*, *FGFR3* and *TP53* (**Figure 6b**). Two of the genes
436 were mutated at least twice as frequently in the normal skin of high-cSCC patients as that of low-
437 cSCC patients: *TP53* (ratio = 3.25) and *FAT1* (ratio = 2.4). Additionally, two less frequently
438 mutated genes, *KRAS* and *HRAS*, were almost exclusively mutated in high-cSCC patients (9 of
439 10). None of these differences reached statistical significance after multiple test correction (data
440 not shown), indicating that larger cohorts will be needed to further explore these potential
441 associations.

442 Although the normal skin of high-cSCC patients contain more mutations per sample,
443 unexpectedly, these mutations were associated with significantly lower VAFs (median=1.0%) than
444 the normal skin of low-cSCC patients (median = 1.3%, $p = 0.011$, Wilcoxon). We found this overall
445 reduction in VAF resulted from a higher number of low-frequency mutations in high-cSCC patients

446 **(Figure 6c)**. For mutations with VAF greater than 1%, the mutations were equally present in high-
447 and low-cSCC patients. However, for low-VAF mutations (defined as <1%), the numbers of
448 mutations per sample were significantly higher in high-cSCC (median = 9.5) than low-cSCC
449 patients (median = 6, $p = 0.032$, Wilcoxon, **Figure 6d**).

450 We next further refined the analysis by focusing on USMs. There were a total of 206 USMs,
451 including 8 CC>TT DNVs. We observed a significantly greater number of USMs in the high-cSCC
452 normal skin samples (median = 11) than the low-cSCC ones (median = 6.5, FDR $p = 0.015$)
453 **(Figure 6e)**. The tumor samples were found to harbor even higher numbers of USMs (median =
454 15.5). The CRCA values, as defined in the primary cohort, were significantly higher in the tumor
455 than the normal skin samples (FDR $p = 0.03$) in the extended cohort. The normal skin samples
456 from high-cSCC patients had slightly higher CRCAs (median = 0.37) than low-cSCC patients
457 (median = 0.31), but the difference was not statistically significant ($p = 0.16$). The CRCA is
458 essentially the sum of VAF values for all detected mutations, normalized for biopsy size. The lack
459 of a significant difference between CRCA values for high-cSCC and low-cSCC skin samples is
460 likely due to the observation that the increased mutations present in high-cSCC samples were
461 enriched for low-frequency mutations (VAF < 1%). We found no significant difference in overall
462 mutation burden, VAF, USMs, or CRCA between normal skin samples collected at 1mm versus
463 6mm from the surgical margin. Lastly, almost all mutations (>99%) were present only in one of
464 two skin samples from the same patient. The absence of shared recurrent mutations across
465 different samples from the same individual indicates that the identified mutations arose
466 independently.

467 **Discussion**

468 Most cancers are initiated by accumulation of somatic mutations^{30,31}. However, early mutations in
469 normal tissues are difficult to detect due to the low abundance and random patterns. Several

470 recent studies demonstrated the feasibility of detecting clonal mutations (CMs) using high-
471 throughput sequencing in various tissue types^{11,12,32}. However, the contribution of these CMs to
472 cancer remains unclear in several ways: how they are generated, what types of mutations are
473 generated by which exogenous and endogenous carcinogens; how the CMs are accumulated
474 and selected by the host microenvironment and inter-clonal competition²⁸; and which mutations
475 contribute or lead to the development of cancer. Indeed, all types of tissues are under the
476 influence of multiple intrinsic and extrinsic factors that vary greatly by individual's lifestyle and
477 environment. Therefore, studying the CMs generated by one specific carcinogen requires
478 comparative studies of matched sample types.

479 To the best of our knowledge, the current study of paired SE and NE skin areas is the first
480 analysis of individual-matched normal human skin to specifically characterize UV radiation's
481 mutational effects. We optimized our detection of UV-induced CMs by: 1) acquiring matched
482 SE/NE skin samples from the same individual to control for aging and other environmental factors
483 unrelated to UV; 2) separating epidermal from dermal layers to decrease non-epidermal
484 background DNA quantity; and 3) ultra-deep DNA sequencing for maximized sensitivity followed
485 by error-suppression to exclude sequencing and alignment errors. Consistent with previous
486 studies^{11,12}, CMs were widespread in epidermal samples. As expected, mutation burden and
487 VAFs were significantly elevated in SE samples. The mutational signatures of the current CMs
488 are consistent with those previously found in skin cancers²⁹, supporting the contribution of the
489 CMs to potential ongoing tumorigenesis. Markedly, our unique approach allowed us to gain
490 several important new insights about epidermal CMs. First, we identified the existence of
491 "mutation-exempt" regions in human genomes. Although mutations frequently occur across most
492 of the sequenced regions in NE skin, presumably due to metabolism and aging related factors,
493 no detectable mutations were found in these mutation-exempt regions. It is unclear whether the
494 absence of mutations in these genomic regions is caused by an active protection or a passive

495 selection mechanism involving altered clone fitness. Interestingly, the “mutation-exempt” property
496 of these regions appears to be altered upon exposure to UV radiation, and these regions become
497 highly mutable. Further studies are warranted to explore how this mechanism is abrogated by UV
498 radiation. Second, USMs were significantly enriched in Glutamate Metabotropic Receptor 3
499 (*GRM3*) in SE skin, which was previously identified as a potential therapeutic target in melanoma
500 ³³, but not reported as a cancer driver in cutaneous SCC. Third, we identified six mutations that
501 were almost exclusively mutated in SE skin. All six mutations had been previously reported in
502 human cutaneous squamous cell carcinomas in the cBioPortal ²⁷. Among these mutations, *TP53*
503 R248W and G245D were highly recurrent with hundreds of occurrences reported in *COSMIC* ³⁴,
504 indicating that the presence of these mutations may be representative of an early phase of
505 carcinogenesis.

506 Consistent with the current finding that UV-exposure results in higher USM burden, and the
507 known knowledge that UV-exposure directly correlates with the risk of cSCC ³⁵, the results of our
508 extended cohort of cSCC patients provided direct evidence that elevated USM burdens are
509 associated with increased burden of cSCC. Presumably, this burden correlates with risk of future
510 cSCC as well. Unexpectedly, we further discovered that most mutational difference between
511 normal skin of high- and low-cSCC patients derived from low-frequency clones (VAF<1%) but not
512 the “expanded” clones (VAF≥1%). It remains unclear why such difference was not seen in the
513 expanded clones. One potential explanation is that the expanded clones might be under more
514 aggressive immune surveillance, as it has been previously reported that the immune system
515 preferably targets larger clones than smaller ones ³⁶. The low-frequency clones, on the other hand,
516 are less actively monitored by the immune system and may more truthfully represent the level of
517 ongoing mutational activity or genomic instability. In any case, the total USM burden in sun-
518 exposed skin of patients with cSCC may be a more accurate measure of skin cancer risk than
519 VAF or clonal area.

520 Our approach was directed by future clinical utilities, focusing on quantitative measurement of
521 UV-induced DNA damage for sun-protection, and cSCC patient risk stratification. These results
522 demonstrate the feasibility of using a small panel of genomic regions (5.5 kb) to quantitatively
523 measure UV-induced CMs. We established Relative Cumulative Clonal Area (CRCA) as a
524 combined measure of mutation burden and relative abundance, which was strongly correlated
525 with sun exposure status, but not with cSCC burden in sun-exposed skin. In the current study, we
526 found the most effective punch size for capturing CMs was 2 mm, which is also clinically favorable
527 as it leaves relatively smaller scars due to the small diameter punch. In future, a non-invasive skin
528 sampling method may provide even wider accessibility to epidermal sampling. In addition, the
529 efficiency of this panel is related to the performance of sequencing method and mutation calling
530 algorithm, which will likely be improved with adoption of more sensitive future methods focusing
531 on the genomic hotspots that are sensitive to UV exposure.

532 The current study focused on the most frequently mutated regions in sun-exposed skin
533 samples defined by the mutations in a previous study ¹². However, we note that many of these
534 regions are mutated in both sun-exposed and non-sun-exposed skin samples, indicating that
535 many mutations in these regions were unrelated to UV exposure. In fact, only 6 of 55 original
536 regions were found to harbor significantly enriched mutations in SE samples. Future studies,
537 including much larger targeted regions, are needed to systematically identify UV-sensitive
538 genomic regions. The skin samples were collected at the same time; therefore, they do not
539 provide longitudinal information about clone initiation and progression. While our analyses of the
540 extended cohort indicate that the burdens of CMs in normal skin are correlated with cancer risk
541 in cSCC patients, this finding needs to be validated in a larger cohort of patients. Future studies
542 including biopsies of both SCC and adjacent normal skin acquired at multiple time points are
543 warranted to unveil the complete role of these CMs in cancer.

544

545 **Conclusions**

546 In summary, this study revealed previously unknown mutational patterns associated with UV-
547 exposure, providing important insights into the early carcinogenic effects of UV radiation. The
548 quantification of CMs has the potential to become a cornerstone for future development of
549 quantitative measures of UV-induced DNA damage, as measured by CRCA, in the clinical setting
550 to monitor early carcinogenesis and highlight the importance of sun protection. The identified
551 association between cSCC burden and mutation status, especially low-frequency CMs, if
552 validated in an expanded cohort, may become a novel biomarker for risk stratification of cSCCs.

553 **Ethics Statements**

554 All specimens in the primary cohort were collected from post-mortem donors collected in
555 collaboration with Buffalo's local organ procurement organization (ConnectLife, formerly Unyts)
556 the Roswell Park's Rapid Tissue Acquisition Program under a Roswell Park approved IRB
557 protocol. Specimens in the expanded cohort were collected from discarded surgical tissue under
558 a Yale University Human Investigation Committee approved protocol.

559 **Availability of data and materials**

560 The datasets used and/or analyzed during the current study are available from the
561 corresponding authors upon request.

562 **Competing interests**

563 None.

564 **Funding**

565 This work was mainly supported by the Roswell Park Alliance Foundation. LW and SL were
566 supported in part by NIH grant U24CA232979. The utilized Genomics and Bioinformatics Shared
567 Resources and Rapid Tissue Acquisition Program at Roswell Park Comprehensive Cancer
568 Center was supported by NCI grant P30CA016056. LW and JX were supported in part by a travel
569 grant from NIH 5U24ES026465. SC was supported by a Career Development Award from the
570 Dermatology Foundation.

571 **Acknowledgments**

572 The authors thank the excellent technical help provided by Paula Pera, MS, assistance with
573 design of the 13-gene custom sequencing panel for the extended cohort provided by Yuemei
574 Zhang, MD (Yale University School of Medicine), and assistance with library generation and
575 sequencing of the extended cohort provided by Mei Zhong, PhD (Yale Stem Cell Center).

576 We dedicate our work to Dr. Oscar Colegio, who passed away suddenly on June 14th, 2020.
577 Oscar was not just a colleague and co-author, he was a passionate and exceptionally empathetic
578 physician, a brilliant researcher, a thoughtful friend. Oscar was among the leading transplant
579 dermatologists in the world. He had a captivating personality and a unique ability to connect ideas
580 and people. This manuscript is a testament to Oscar's ability to bring people together. As we will
581 continue our collaboration Dr. Colegio's insight, mentorship, wit and hard work will be greatly
582 missed.

583 **List of abbreviations**

584 UV – Ultraviolet

585 CM – Clonal mutation

586 NMSC – Nonmelanoma skin cancer

587 SE – Sun-exposed
588 NE – Non-sun-exposed
589 USM – UV-signature mutation
590 NUSM – Non-UV-signature mutation
591 CRCA – Cumulative Relative Clonal Area
592 cSCC – Cutaneous squamous cell carcinoma
593 AK – Actinic keratosis
594 SNV – Single nucleotide variant
595 Indels – Insertions/deletions
596 DNV – Dinucleotide variant
597 CSNV – Cluster of single nucleotide variant
598 MAC – Multi-Nucleotide Variant Annotation Corrector
599 VAF – Variant allele frequency
600

601 **Figure Legends**

602 **Figure 1. Region-specific enrichment of somatic mutations in sun-exposed skin.** a). Graph
603 shows the number of mutations identified within each 100-bp genomic target window grouped by
604 SE and NE skin types. b). The overall gene-level number of mutations from SE and NE samples.
605 Stars indicate the segments or genes where mutations are significantly enriched in the SE
606 samples (FDR p values: *** $p < 0.001$; + $p < 0.1$).

607 **Figure 2. Hotspots and mutations associated with UV-exposure.** a). All mutations are ordered
608 by their genomic locations. X-axis: the order of the mutation's genomic location. Y-axis: variant
609 allele fraction (VAF) of individual mutations. Color depicts the gene harboring the mutations. The
610 three genes demonstrating significant difference between SE and NE, either on the gene level or
611 segment level, were labeled on top (*TP53*, *GRM3*, *NOTCH1*). One specific mutation with elevated
612 VAFs (*NOTCH1* E424K) is indicated with a red arrow. b). The VAF of the six individual mutations
613 that are significantly enriched in SE vs NE epidermis in the primary discovery (green) and
614 validation (brown) data sets. The dotted red line represents median VAF of all mutations and
615 black lines indicate the median of each group. c). The predicted protein complex structure of
616 *NOTCH1* and *DLL4* to show the position of the mutant E424K and the interacting partners, *DLL4*
617 K189/R191, in wild type.

618 **Figure 3. UV-induced DNA damage assessed by USMs.** a). Only UV-signature mutations are
619 associated with sun-exposure status. Left: higher numbers of USMs are found in SE than NE skin.
620 Right: NUSMs are almost equally presented in SE and NE samples. Red dotted line indicates
621 high-VAF (> 0.1). Black dotted circle indicates extra USMs in SE skin compared with NE skin. b).
622 The numbers of mutations by each amino-acid-change type found in SE and NE skin, grouped
623 by USMs and NUSMs. Overall distribution c.) of the numbers of USMs per sample and d.) the
624 VAFs of the mutations using the 2 mm punch size. Inside the violin plots: black dots - original

625 data points from individual samples; yellow dot with bar - averaged value with standard deviation.
626 SE samples are associated with higher numbers of USMs, as well as higher VAFs indicating
627 potential larger clones. e). Cumulative Relative Clonal Area (CRCA) was developed to represent
628 the overall percentage of the biopsied skin areas that are covered by clonal mutations. In all 13
629 current individuals, CRCAs were higher in SE than in the matched NE group, with the ratios of
630 SE/NE ranged from 1.4 to 55.0 (mean = 11.2). Statistical tests used: figure 3b, Fisher's exact test
631 with multiple test correction implemented using the FDR method; figure 3c, 3d: Wilcoxon test;
632 * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

633 **Figure 4. Optimization of punch size for detecting USMs.** a). A representative figure showing
634 one representative punch of each collection size. We selected the sample with the highest number
635 of mutations under each size for easy illustration. Every mutation is plotted as a dot with its size
636 calculated to match the clonal area harboring the mutation. One punch size, 3 mm was not shown
637 as it was obtained by cutting a 6 mm punch into quarters. b). In the discovery cohort, 2 mm was
638 found to be the most efficient size in differentiating CRCA from SE and NE skin samples by p
639 value. c). Distribution of numbers of USMs per sample at each punch size, after combining both
640 the discovery and validation cohorts. d). VAF of USM detected in different size punches. The size
641 of the dot indicates the approximate relative area of cells containing the mutation. In SE samples,
642 VAFs of USMs detected from larger punches are associated with smaller variations.

643 **Figure 5. Mutational contexts associated with UV-exposure.** a). Each dot represents a specific
644 mutation context of SNVs and DNVs. X-axis: the total numbers of mutations of each context; y-
645 axis: p value of the context for differentiating SE and NE skin, shown as $-\log(p)$. The dotted line
646 indicate $p < 0.05$ (the above area). None of the NUSM contexts was significant. b). Further
647 refinement of USM contexts by depicting the numbers of mutations in SE and NE skin for all
648 current USM contexts. Mutation contexts are ordered by the p value of SE vs NE in an increasing

649 order from left to right. Multiple test correction was implemented using the FDR method. The
650 dotted line indicates FDR $p < 0.05$ (the left side).

651

652 **Figure 6.**

653 **Clonal mutations are correlated with cSCC burden.** a). Violin plots depicting the overall
654 distribution of somatic mutations in each sample, ordered by sample type. b). Mutation numbers
655 by genes in the normal skin. NS (high) - normal skin from high-cSCC patients; NS(low) - normal
656 skin from low-cSCC patients. c). High-cSCC patients are associated with increased low-VAF
657 ($< 1\%$) mutations. Histogram depicting the distribution of VAFs of the detected mutations in normal
658 skin separated by cSCC burden. The dotted oval highlights the increased low-VAF mutations in
659 the normal skin of high-cSCC patients compared with low-cSCC patients. d) Number of mutations
660 per sample in normal skin, separated by high- ($\geq 1\%$) and low- ($< 1\%$) VAFs; e) Number of USMs
661 per sample in high- and low- cSCC normal skin (NS), and cSCC tumors. Shape indicates the two
662 normal skin samples from each patient, taken either 1mm (circle) or 6mm (triangle) from the
663 surgical margin.

664 **Tables**

665 Table 1. Patient and sample cohort

Patient	Control (dermis)	Epidermis SE/NE pairs						Total SE/NE pairs
		1 mm	2 mm	3 mm*	4 mm	6 mm	2 mm [#]	
Pt1	1		5	4	1	1		11
Pt2	1	3	5	4	1	1		14
Pt3	1	3	5	4	1	1		14
Pt4	1	3	3					6
Pt5	1	3	3	3	3	3	5	20
Pt6	1	3	3	3	3	3	5	20
Pt7	1	3	3	3	3	3	5	20
Pt8	1	3	3	3	3	3	5	20
Pt9	1	3	3	3	3	3	5	20
Pt10	1	3	3	3	3	3	5	20
Pt11	1	3	3	3	3	3	5	20
Pt12	1	3	3	3	3	3	5	20
Pt13	2	3	3	3	3	3	5	20
Total	14	36	45	39	30	30	45	225

666 * 3 mm punches were obtained by cutting 6 mm punches into quarters

667 # Validation cohort containing only 2 mm punches

668

669 **References**

- 670 1 Rogers, H. W., Weinstock, M. A., Feldman, S. R. & Coldiron, B. M. Incidence Estimate of
671 Nonmelanoma Skin Cancer (Keratinocyte Carcinomas) in the U.S. Population, 2012.
672 *JAMA Dermatol* **151**, 1081-1086, doi:10.1001/jamadermatol.2015.1187 (2015).
- 673 2 Koh, H. K., Geller, A. C., Miller, D. R., Grossbart, T. A. & Lew, R. A. Prevention and early
674 detection strategies for melanoma and skin cancer. Current status. *Arch Dermatol* **132**,
675 436-443 (1996).
- 676 3 Martejn, J. A., Lans, H., Vermeulen, W. & Hoeijmakers, J. H. Understanding nucleotide
677 excision repair and its roles in cancer and ageing. *Nature reviews. Molecular cell biology*
678 **15**, 465-481, doi:10.1038/nrm3822 (2014).
- 679 4 Brash, D. E. UV signature mutations. *Photochem Photobiol* **91**, 15-26,
680 doi:10.1111/php.12377 (2015).
- 681 5 Wikonkal, N. M. & Brash, D. E. Ultraviolet radiation induced signature mutations in
682 photocarcinogenesis. *J Investig Dermatol Symp Proc* **4**, 6-10 (1999).
- 683 6 Marks, R., Rennie, G. & Selwood, T. S. Malignant transformation of solar keratoses to
684 squamous cell carcinoma. *Lancet* **1**, 795-797, doi:10.1016/s0140-6736(88)91658-3
685 (1988).
- 686 7 Ling, G. *et al.* Persistent p53 mutations in single cells from normal human skin. *Am J*
687 *Pathol* **159**, 1247-1253, doi:10.1016/S0002-9440(10)62511-4 (2001).
- 688 8 Brash, D. E. Cancer. Preprocancer. *Science* **348**, 867-868, doi:10.1126/science.aac4435
689 (2015).
- 690 9 Urano, Y. *et al.* Frequent p53 accumulation in the chronically sun-exposed epidermis and
691 clonal expansion of p53 mutant cells in the epidermis adjacent to basal cell carcinoma. *J*
692 *Invest Dermatol* **104**, 928-932 (1995).
- 693 10 Williams, C. *et al.* Clones of normal keratinocytes and a variety of simultaneously present
694 epidermal neoplastic lesions contain a multitude of p53 gene mutations in a xeroderma
695 pigmentosum patient. *Cancer Res* **58**, 2449-2455 (1998).
- 696 11 Yizhak, K. *et al.* RNA sequence analysis reveals macroscopic somatic clonal expansion
697 across normal tissues. *Science* **364**, doi:10.1126/science.aaw0726 (2019).
- 698 12 Martincorena, I. *et al.* Tumor evolution. High burden and pervasive positive selection of
699 somatic mutations in normal human skin. *Science* **348**, 880-886,
700 doi:10.1126/science.aaa6806 (2015).
- 701 13 Wei, L. *et al.* Accurate Quantification of Residual Cancer Cells in Pelvic Washing Reveals
702 Association with Cancer Recurrence Following Robot-Assisted Radical Cystectomy. *J*
703 *Urol* **201**, 1105-1114, doi:10.1097/JU.000000000000142 (2019).
- 704 14 Wei, L. *et al.* Pitfalls of improperly procured adjacent non-neoplastic tissue for somatic
705 mutation analysis using next-generation sequencing. *BMC medical genomics* **9**, 64,
706 doi:10.1186/s12920-016-0226-1 (2016).
- 707 15 Tang, J. *et al.* The genomic landscapes of individual melanocytes from human skin.
708 *bioRxiv*, 2020.2003.2001.971820, doi:10.1101/2020.03.01.971820 (2020).

- 709 16 Huss, W. J. *et al.* Comparison of SureSelect and Nextera Exome Capture Performance in
710 Single-Cell Sequencing. *Human heredity* **83**, 153-162, doi:10.1159/000490506 (2018).
- 711 17 Gamble, R. G. *et al.* Sun damage in ultraviolet photographs correlates with phenotypic
712 melanoma risk factors in 12-year-old children. *J Am Acad Dermatol* **67**, 587-597,
713 doi:10.1016/j.jaad.2011.11.922 (2012).
- 714 18 Creidi, P. *et al.* Profilometric evaluation of photodamage after topical retinaldehyde and
715 retinoic acid treatment. *J Am Acad Dermatol* **39**, 960-965 (1998).
- 716 19 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
717 transform. *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).
- 718 20 Liu, Q. *et al.* SeqSQC: A Bioconductor Package for Evaluating the Sample Quality of Next-
719 generation Sequencing Data. *Genomics Proteomics Bioinformatics* **17**, 211-218,
720 doi:10.1016/j.gpb.2018.07.006 (2019).
- 721 21 Saunders, C. T. *et al.* Strelka: accurate somatic small-variant calling from sequenced
722 tumor-normal sample pairs. *Bioinformatics* **28**, 1811-1817,
723 doi:10.1093/bioinformatics/bts271 (2012).
- 724 22 Wei, L. *et al.* MAC: identifying and correcting annotation for multi-nucleotide variations.
725 *BMC genomics* **16**, 569, doi:10.1186/s12864-015-1779-7 (2015).
- 726 23 Newman, A. M. *et al.* Integrated digital error suppression for improved detection of
727 circulating tumor DNA. *Nat Biotechnol* **34**, 547-555, doi:10.1038/nbt.3520 (2016).
- 728 24 Luca, V. C. *et al.* Structural biology. Structural basis for Notch1 engagement of Delta-like
729 4. *Science* **347**, 847-853, doi:10.1126/science.1261093 (2015).
- 730 25 Blanpain, C., Lowry, W. E., Pasolli, H. A. & Fuchs, E. Canonical notch signaling functions
731 as a commitment switch in the epidermal lineage. *Genes Dev* **20**, 3022-3035,
732 doi:10.1101/gad.1477606 (2006).
- 733 26 Lefort, K. & Dotto, G. P. Notch signaling in the integrated control of keratinocyte
734 growth/differentiation and tumor suppression. *Semin Cancer Biol* **14**, 374-386,
735 doi:10.1016/j.semcancer.2004.04.017 (2004).
- 736 27 Gao, J. *et al.* Integrative analysis of complex cancer genomics and clinical profiles using
737 the cBioPortal. *Science signaling* **6**, pl1, doi:10.1126/scisignal.2004088 (2013).
- 738 28 Colom, B. *et al.* Spatial competition shapes the dynamic mutational landscape of normal
739 esophageal epithelium. *Nat Genet*, doi:10.1038/s41588-020-0624-3 (2020).
- 740 29 Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**,
741 415-421, doi:10.1038/nature12477 (2013).
- 742 30 Wu, S., Powers, S., Zhu, W. & Hannun, Y. A. Substantial contribution of extrinsic risk
743 factors to cancer development. *Nature* **529**, 43-47, doi:10.1038/nature16166 (2016).
- 744 31 Tomasetti, C. & Vogelstein, B. Cancer etiology. Variation in cancer risk among tissues can
745 be explained by the number of stem cell divisions. *Science* **347**, 78-81,
746 doi:10.1126/science.1260825 (2015).
- 747 32 Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus with age.
748 *Science*, doi:10.1126/science.aau3879 (2018).
- 749 33 Kunz, M. The genetic basis of new treatment modalities in melanoma. *Curr Drug Targets*
750 **16**, 233-248, doi:10.2174/1389450116666150204112138 (2015).

- 751 34 Forbes, S. A. *et al.* The Catalogue of Somatic Mutations in Cancer (COSMIC). *Curr Protoc*
752 *Hum Genet* **Chapter 10**, Unit 10 11, doi:10.1002/0471142905.hg1011s57 (2008).
- 753 35 Johnson, T. M., Rowe, D. E., Nelson, B. R. & Swanson, N. A. Squamous cell carcinoma
754 of the skin (excluding lip and oral mucosa). *J Am Acad Dermatol* **26**, 467-484,
755 doi:10.1016/0190-9622(92)70074-p (1992).
- 756 36 Gejman, R. S. *et al.* Rejection of immunogenic tumor clones is limited by clonal fraction.
757 *Elife* **7**, doi:10.7554/eLife.41090 (2018).

758

Figure 1. Region-specific enrichment of somatic mutations in sun-exposed skin

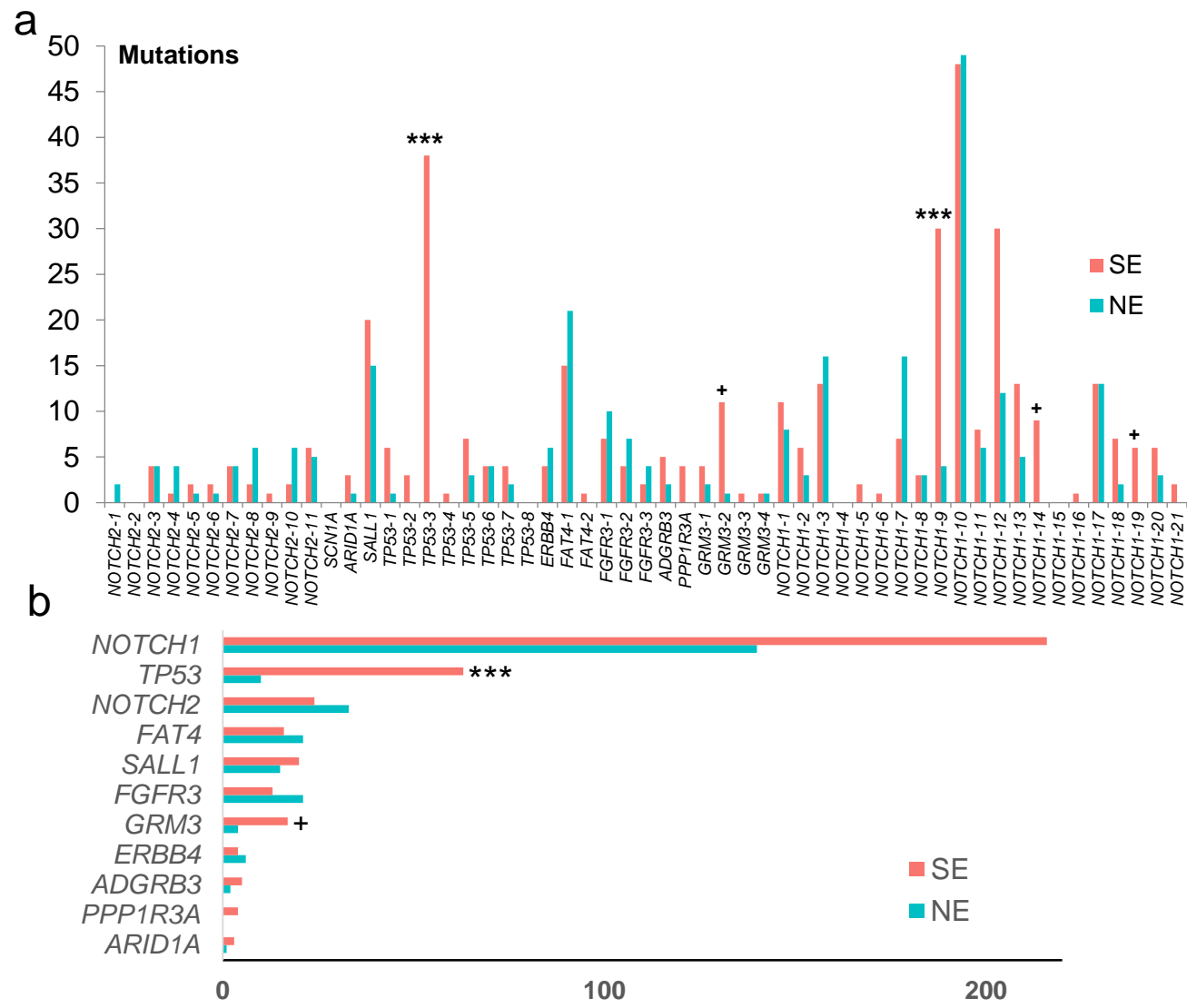
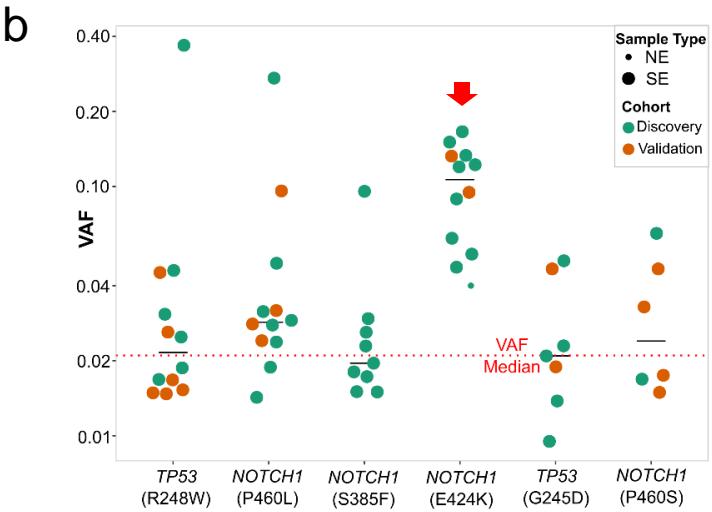
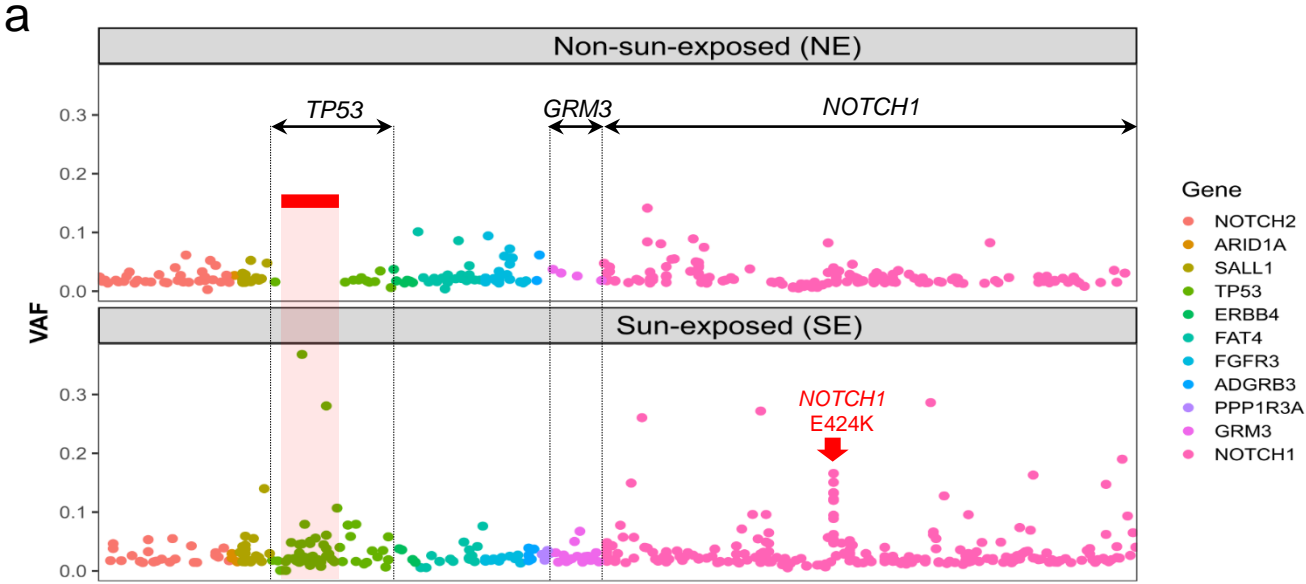


Figure 2. Hotspots and mutations associated with UV-exposure.



c

NOTCH1

EGF11^{424(K)}

DSL

DLL4

Figure 3. UV-induced DNA damage assessed by USMs

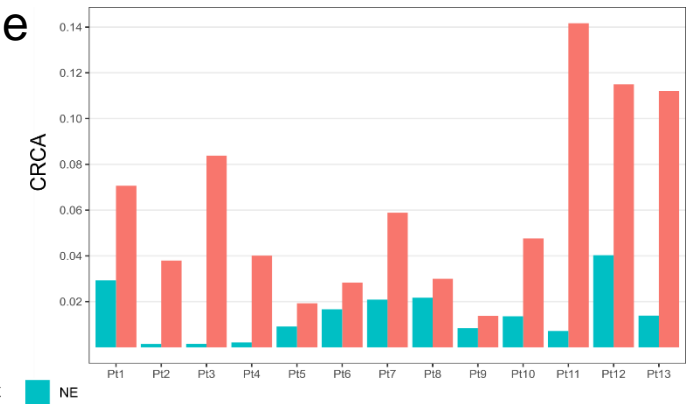
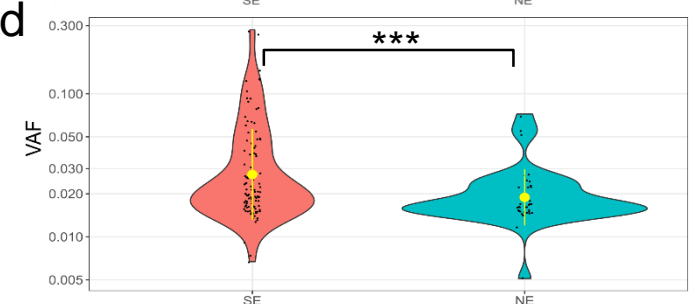
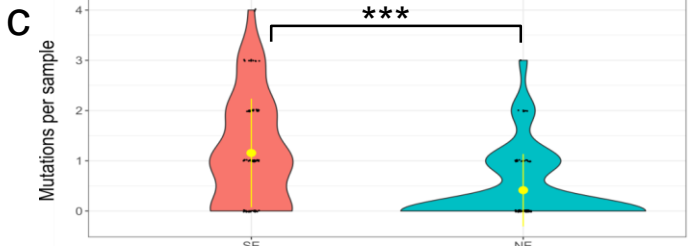
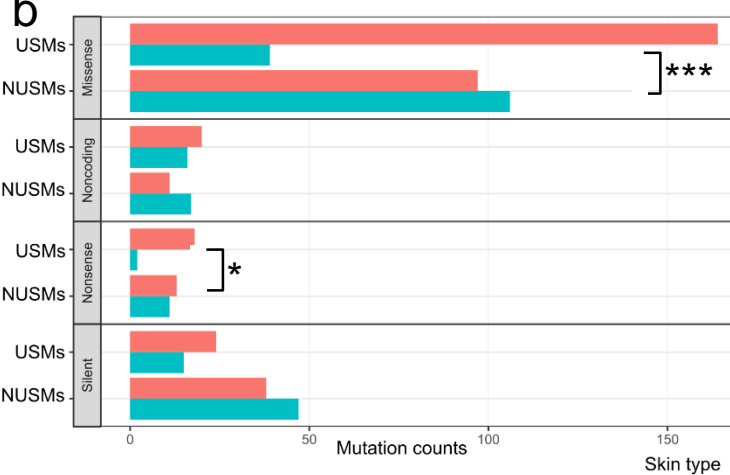
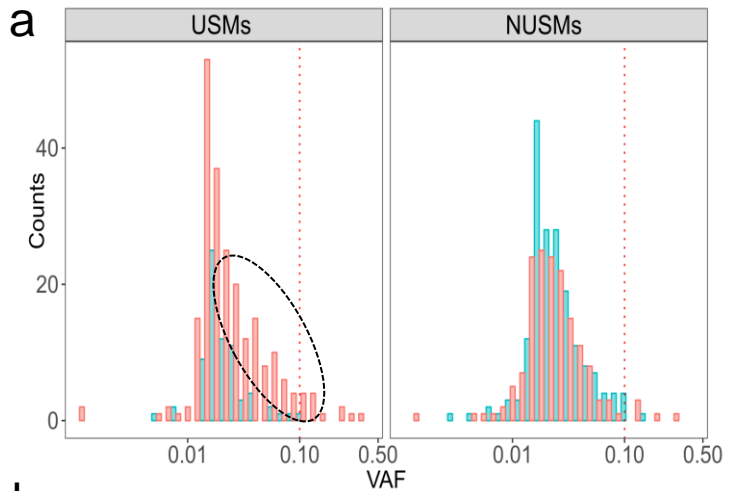


Figure 4. Optimization of punch size for detecting USMs

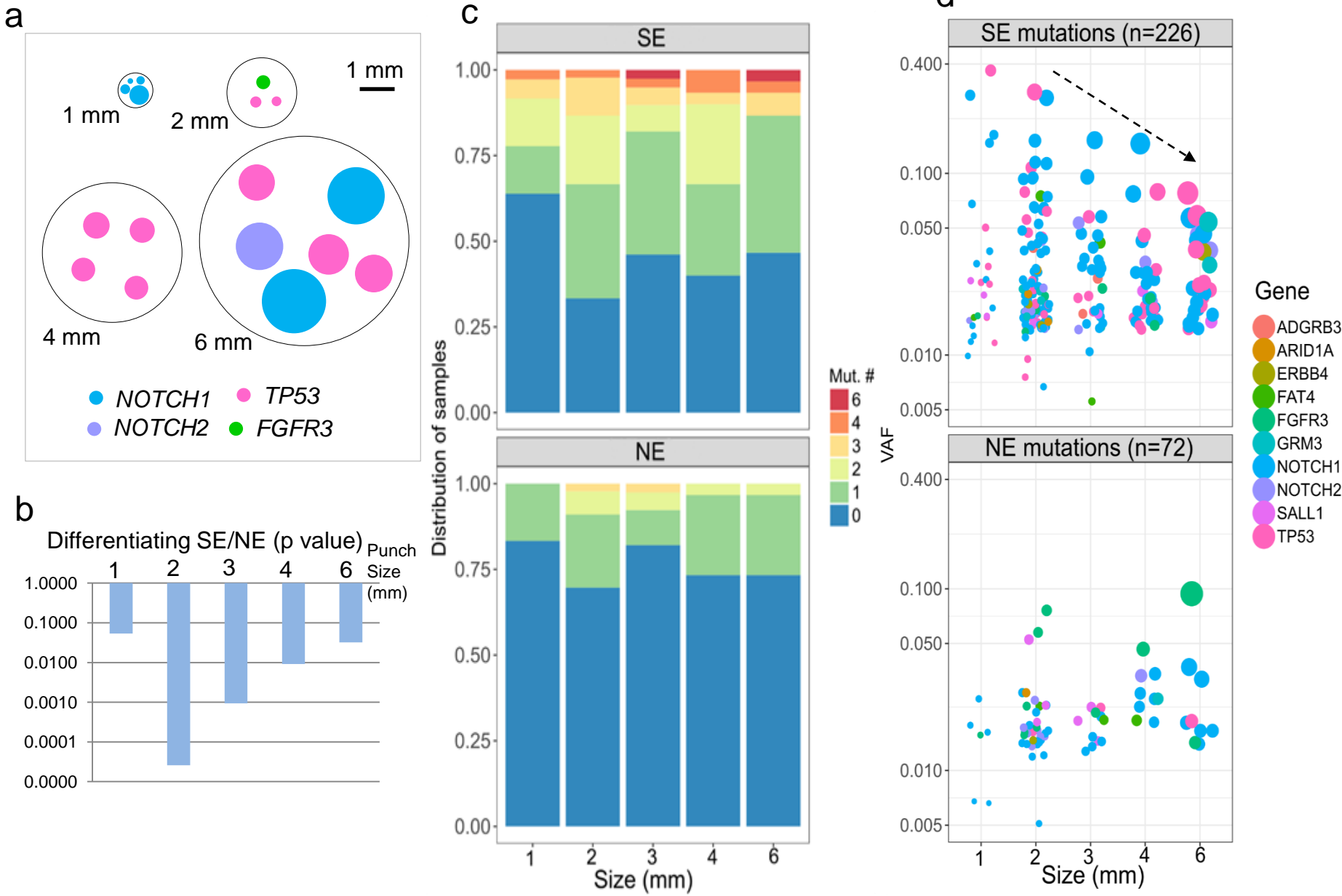


Figure 5. Mutational contexts associated with UV-exposure

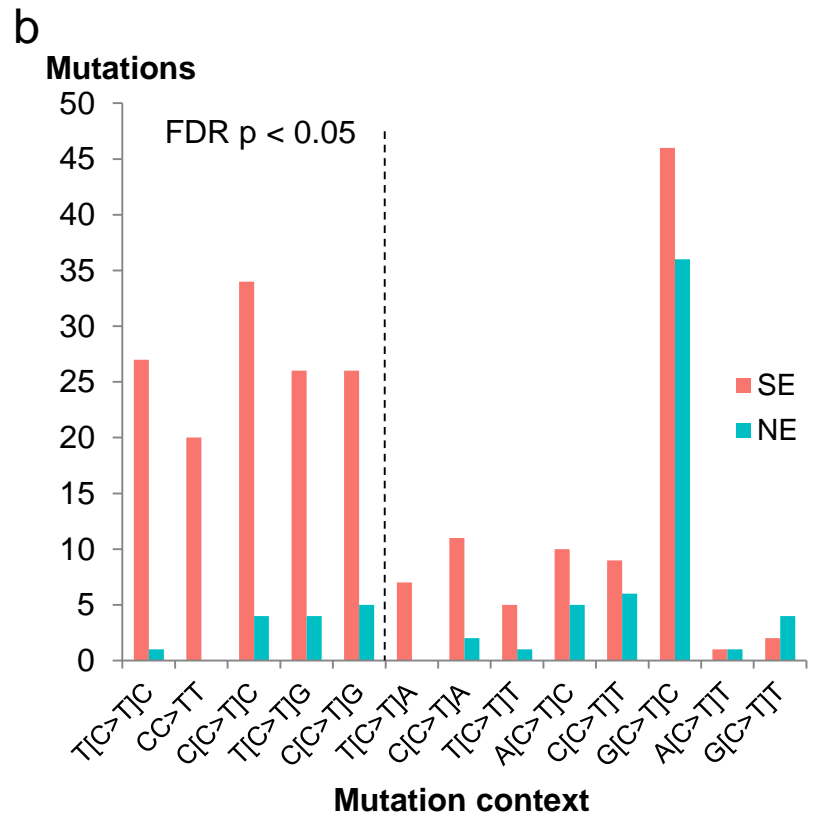
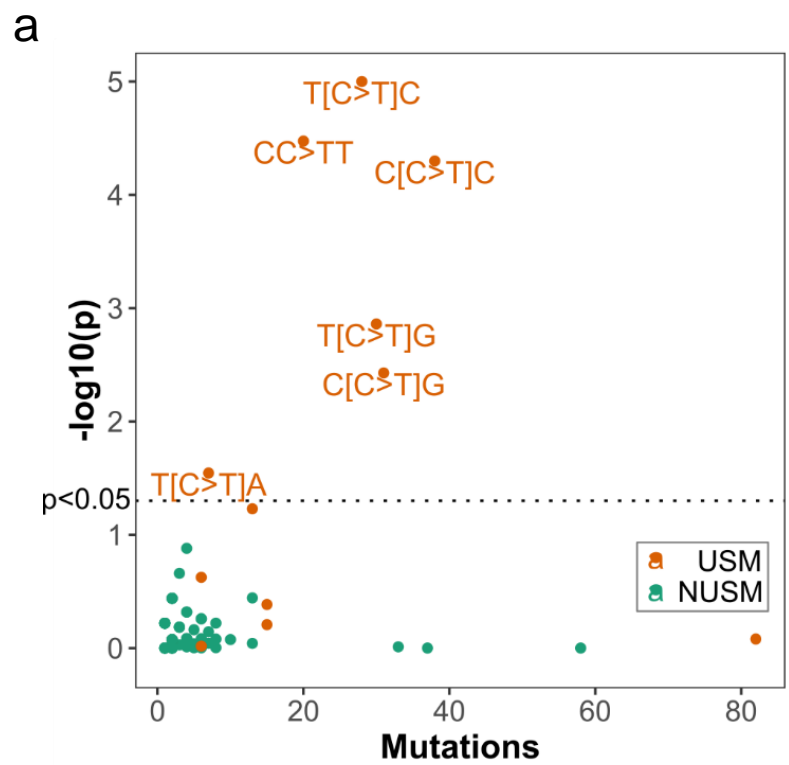


Figure 6. Clonal mutation burden correlates with skin cancer risk

