

1 **Geographic variation of mutagenic exposures in kidney cancer genomes**

2

3 Sergey Senkin¹, Sarah Moody², Marcos Díaz-Gay^{3,4,5}, Behnoush Abedi-Ardekani¹, Thomas

4 Cattiaux¹, Aida Ferreira-Iglesias¹, Jingwei Wang², Stephen Fitzgerald², Mariya

5 Kazachkova^{3,6,5}, Raviteja Vangara^{3,4,5}, Anh Phuong Le², Erik N. Bergstrom^{3,4,5}, Azhar

6 Khandekar^{3,4,5}, Burçak Otlu^{3,4,5,7}, Saamin Cheema², Calli Latimer², Emily Thomas², Joshua

7 Ronald Atkins⁸, Karl Smith-Byrne⁸, Ricardo Cortez Cardoso Penha¹, Christine Carreira⁹,

8 Priscilia Chopard¹, Valérie Gaborieau¹, Pekka Keski-Rahkonen¹⁰, David Jones², Jon W.

9 Teague², Sophie Ferlicot¹¹, Mojgan Asgari¹², Surasak Sangkhathat¹³, Worapat

10 Attawettayanon¹⁴, Beata Świątkowska¹⁵, Sonata Jarmalaite^{16,17}, Rasa Sabaliauskaite¹⁶,

11 Tatsuhiro Shibata^{18,19}, Akihiko Fukagawa^{19,20}, Dana Mates²¹, Viorel Jinga²², Stefan Rascu²²,

12 Mirjana Mijuskovic²³, Slavisa Savic²⁴, Sasa Milosavljevic²⁵, John M.S. Bartlett²⁶, Monique

13 Albert²⁷, Larry Phouthavongsy²⁸, Patricia Ashton-Prolla^{29,30}, Mariana R. Botton³¹, Brasil Silva

14 Neto^{32,33}, Stephania Martins Bezerra³⁴, Maria Paula Curado³⁵, Stênio de Cássio

15 Zequi^{36,37,38,39}, Rui Manuel Reis^{40,41}, Eliney Faria⁴², Nei Soares Menezes⁴³, Renata Spagnoli

16 Ferrari⁴², Rosamonde E. Banks⁴⁴, Naveen S. Vasudev⁴⁴, David Zaridze⁴⁵, Anush Mukeriya⁴⁵,

17 Oxana Shangina⁴⁵, Vsevolod Matveev⁴⁶, Lenka Foretova⁴⁷, Marie Navratilova⁴⁷, Ivana

18 Holcatova^{48,49}, Anna Hornakova⁵⁰, Vladimir Janout⁵¹, Mark Purdue⁵², Stephen J. Chanock⁵²,

19 Per Magne Ueland⁵³, Mattias Johansson¹, James McKay¹, Ghislaine Scelo⁵⁴, Estelle

20 Chanudet⁵⁵, Laura Humphreys², Ana Carolina de Carvalho¹, Sandra Perdomo¹, Ludmil B.

21 Alexandrov^{3,4,5}, Michael R. Stratton², Paul Brennan^{1*}

22

23 ¹Genomic Epidemiology Branch, International Agency for Research on Cancer (IARC/WHO),

24 Lyon, France, ²Cancer, Ageing and Somatic Mutation, Wellcome Sanger Institute,

25 Cambridge, UK, ³Department of Cellular and Molecular Medicine, University of California

26 San Diego, La Jolla, USA, ⁴Department of Bioengineering, University of California San

27 Diego, La Jolla, USA, ⁵Moore's Cancer Center, University of California San Diego, La Jolla,

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

28 USA, ⁶Biomedical Sciences Graduate Program, University of California San Diego, La Jolla,
29 USA, ⁷Department of Health Informatics, Graduate School of Informatics, Middle East
30 Technical University, Ankara, Turkey, ⁸Cancer Epidemiology Unit, The Nuffield Department
31 of Population Health, University of Oxford, Oxford, UK, ⁹Evidence Synthesis and
32 Classification Branch, International Agency for Research on Cancer (IARC/WHO), Lyon,
33 France, ¹⁰Nutrition and Metabolism Branch, International Agency for Research on Cancer
34 (IARC/WHO), Lyon, France, ¹¹Service d'Anatomie Pathologique, Assistance Publique-
35 Hôpitaux de Paris, Univeristé Paris-Saclay, Le Kremlin-Bicêtre, France, ¹²Oncopathology
36 Research Center, Iran University of Medical Sciences, Tehran, Iran, ¹³Translational Medicine
37 Research Center, Faculty of Medicine, Prince of Songkla University, Hat Yai, Thailand,
38 ¹⁴Department of Surgery, Urology, Faculty of Medicine, Prince of Songkla University, Hat Yai,
39 Thailand, ¹⁵Department of Environmental Epidemiology, Nofer Institute of Occupational
40 Medicine, Łódź, Poland, ¹⁶Laboratory of Genetic Diagnostic, National Cancer Institute,
41 Vilnius, Lithuania, ¹⁷Department of Botany and Genetics, Institute of Biosciences, Vilnius
42 University, Vilnius, Lithuania, ¹⁸Laboratory of Molecular Medicine, The Institute of Medical
43 Science, The University of Tokyo, Minato-ku, Japan, ¹⁹Division of Cancer Genomics,
44 National Cancer Center Research Institute, Chuo-ku, Japan, ²⁰Department of Pathology,
45 Graduate School of Medicine, The University of Tokyo, Bunkyo-ku, Japan, ²¹Occupational
46 Health and Toxicology, National Center for Environmental Risk Monitoring, National Institute
47 of Public Health, Bucharest, Romania, ²²Urology Department, "Carol Davila" University of
48 Medicine and Pharmacy - "Prof. Dr. Th. Burgele" Clinical Hospital, Bucharest, Romania,
49 ²³Clinic of Nefrology, Faculty of Medicine, Military Medical Academy, Belgrade, Serbia,
50 ²⁴Department of Urology, University Hospital "Dr D. Misovic" Clinical Center, Belgrade,
51 Serbia, ²⁵International Organization for Cancer Prevention and Research, Belgrade, Serbia,
52 ²⁶Cancer Research UK Edinburgh Centre, Institute of Genetics and Cancer, University of
53 Edinburgh, Edinburgh, Scotland, ²⁷Centre for Biodiversity Genomics, University of Guelph,
54 Guelph, Canada, ²⁸Ontario Tumour Bank, Ontario Institute for Cancer Research, Toronto,
55 Canada, ²⁹Experimental Research Center, Genomic Medicine Laboratory, Hospital de

56 Clínicas de Porto Alegre, Porto Alegre, Brazil, ³⁰Post-Graduate Program in Genetics and
57 Molecular Biology, Universidade Federal do Rio Grande do Sul, Porto Alegre, Brazil,
58 ³¹Diagnostic Laboratory Service, Personalized Medicine, Hospital de Clínicas de Porto
59 Alegre, Porto Alegre, Brazil, ³²Service of Urology, Hospital de Clínicas de Porto Alegre, Porto
60 Alegre, Brazil, ³³Post-Graduate Program in Medicine: Surgical Sciences, Universidade
61 Federal do Rio Grande do Sul, Porto Alegre, Brazil, ³⁴Department of Anatomic Pathology,
62 A.C. Camargo Cancer Center, São Paulo, Brazil, ³⁵Department of Epidemiology, A.C.
63 Camargo Cancer Center, São Paulo, Brazil, ³⁶Department of Urology, A.C. Camargo Cancer
64 Center, São Paulo, Brazil, ³⁷National Institute for Science and Technology in Oncogenomics
65 and Therapeutic Innovation, A.C. Camargo Cancer Center, São Paulo, Brazil, ³⁸Latin
66 American Renal Cancer Group – LARCG, São Paulo, Brazil, ³⁹Department of Surgery,
67 Division of Urology, Sao Paulo Federal University - UNIFESP, São Paulo, Brazil, ⁴⁰Molecular
68 Oncology Research Center, Barretos Cancer Hospital, Brazil, ⁴¹Life and Health Sciences
69 Research Institute (ICVS), School of Medicine, Minho University, Braga, Portugal,
70 ⁴²Department of Urology, Barretos Cancer Hospital, Brazil, ⁴³Department of Pathology,
71 Barretos Cancer Hospital, Brazil, ⁴⁴Leeds Institute of Medical Research at St James's,
72 University of Leeds, Leeds, UK, ⁴⁵Clinical Epidemiology, N.N.Blokhin National Medical
73 Research Centre of Oncology, Moscow, Russia, ⁴⁶Department of Urology, N.N.Blokhin
74 National Medical Research Centre of Oncology, Moscow, Russia, ⁴⁷Department of Cancer
75 Epidemiology and Genetics, Masaryk Memorial Cancer Institute, Brno, Czech Republic,
76 ⁴⁸Institute of Public Health & Preventive Medicine, 2nd Faculty of Medicine, Charles
77 University, Prague, Czech Republic, ⁴⁹Department of Oncology, 2nd Faculty of Medicine,
78 Charles University and Motol University Hospital, Prague, Czech Republic, ⁵⁰Institute of
79 Hygiene & Epidemiology, 1st Faculty of Medicine, Charles University, Prague, Czech
80 Republic, ⁵¹Faculty of Health Sciences, Palacky University, Olomouc, Czech Republic,
81 ⁵²Division of Cancer Epidemiology and Genetics, National Cancer Institute, Rockville, USA,
82 ⁵³Bevital AS, Bergen, Norway, ⁵⁴Observational & Pragmatic Research Institute Pte Ltd,

83 Singapore, Singapore, ⁵⁵Department of Pathology, Radboud University Medical Centre,
84 Nijmegen, Netherlands

85

86 These authors contributed equally: Sergey Senkin, Sarah Moody

87 * Corresponding author: Paul Brennan

88

89 **ABSTRACT**

90 International differences in the incidence of many cancer types indicate the existence of
91 carcinogen exposures which make a substantial contribution to cancer burden, vary
92 geographically, and have underlying agents thus far unidentified by conventional
93 epidemiology¹. This pertains to clear cell renal cell carcinomas (ccRCC), for which obesity,
94 hypertension, and tobacco smoking are risk factors but do not explain its geographical variation
95 in incidence². Some carcinogens generate somatic mutations and past exposures can be
96 inferred from the patterns of mutations found in cancer genomes. Therefore, we sequenced
97 the whole genomes of 962 ccRCC from 11 countries of varying incidence. Somatic mutation
98 profiles differed between countries. In Romania, Serbia and Thailand, mutational signatures
99 likely caused by extracts of Aristolochia plants were present in most cases and rare elsewhere.
100 In Japan, a mutational signature of unknown cause was found in >70% cases and <2%
101 elsewhere. Another mutational signature of unknown cause was ubiquitous and associated
102 with kidney cancer incidence rates ($p\text{-value} < 6 \times 10^{-18}$), with higher numbers of mutations in
103 countries with higher risk. Known signatures of tobacco smoking correlated with tobacco
104 consumption, but no signature was associated with obesity or hypertension suggesting non-
105 mutagenic mechanisms of action underlying these risk factors. The results indicate the
106 existence of multiple, widespread, geographically variable mutagenic exposures to known and
107 unknown agents, which may contribute to the incidence of kidney cancer.

108

109

110 INTRODUCTION

111 The incidence rates of most adult cancers vary substantially between geographical regions
112 and many such differences are unexplained by known risk factors¹. Together with unexplained
113 trends in incidence over time, this indicates the likely presence of unknown environmental or
114 lifestyle causes for many cancer types¹. Kidney cancer, for example, has particularly high
115 incidence rates in Central and Northern Europe, notably in the Czech Republic and Lithuania,
116 and has shown increasing incidence in high income countries in recent decades (**Fig. 1**)². Most
117 kidney cancers are clear cell renal cell carcinomas (ccRCC)³ for which obesity, hypertension
118 and tobacco smoking are known risk factors². However, these account for <50% of the global
119 ccRCC burden and do not explain geographical or temporal incidence trends. Recently,
120 evidence has also emerged of increased risk associated with environmental exposure to per-
121 and polyfluoroalkyl substances (PFAS)³, industrial chemicals used in a wide range of
122 consumer and industrial products.

123

124 Characterization of mutational signatures within cancer genomes⁴ is an approach
125 complementary to conventional epidemiology for investigating unknown causes of cancer.
126 Most cancers contain thousands of somatic mutations that have occurred over the lifetime of
127 the individual. These can be caused by endogenous cellular processes, such as imperfect
128 DNA replication and repair, or by exposure to exogenous environmental or lifestyle mutagens
129 such as ultraviolet radiation in sunlight and compounds in cigarette smoke. Mutational
130 signatures are the patterns of somatic mutation imprinted on genomes by individual mutational
131 processes. Analysis of thousands of cancer genome sequences from most cancer types has
132 established a set of reference mutational signatures including 71 single base substitution
133 (SBS) or doublet base substitution (DBS) signatures, and 18 small insertion and deletion (ID)
134 signatures⁵. A possible etiology has been suggested for 47 SBS/DBS signatures and nine ID
135 signatures.

136

137 Previous ccRCC genome sequencing studies have included relatively modest numbers of
138 individuals from a small number of countries with limited variation in ccRCC incidence^{6–10} and
139 have not comprehensively examined associations between ccRCC risk factors and mutational
140 signatures. To detect the activity of unknown carcinogens involved in ccRCC development and
141 to investigate the mechanisms of action of known risk factors, we generated and analyzed
142 epidemiological and whole genome sequencing data from a large international series of
143 ccRCC¹¹.

144

145 **RESULTS**

146 A total of 962 ccRCC cases from 11 countries in four continents were studied, encompassing:
147 Czech Republic ($n=259$), Russia ($n=216$), United Kingdom ($n=115$), Brazil ($n=96$), Canada
148 ($n=73$), Serbia ($n=69$), Romania ($n=64$), Japan ($n=36$), Lithuania ($n=16$), Poland ($n=13$), and
149 Thailand ($n=5$; **Fig. 1**; **Table 1**; **Methods**). These encompass a broad range of ccRCC
150 incidence, from the highest global age-standardized rates (ASRs) of Lithuania and Czech
151 Republic (ASRs of 14.5 and 14.4/100,000 respectively) to the relatively low rates of Brazil and
152 Thailand (ASRs of 4.5 and 1.8/100,000 respectively)¹². Epidemiological questionnaire data
153 were available on sex, age at diagnosis, and important risk factors including body mass index
154 (BMI), hypertension, and tobacco smoking (**Table 1**). DNAs from ccRCCs and blood from the
155 same individuals were extracted and whole genome sequenced to average coverage of 54-
156 fold and 31-fold, respectively.

157

158 Somatic mutation burdens in the 962 ccRCC genomes ranged from 803 to 45,376 (median
159 5,093) for single base substitutions (SBS), 2 to 240 (median 53) for doublet base substitutions
160 (DBS), and 10 to 14,770 (median 695) for small insertions and deletions (**Supplementary**
161 **Table 1**). The average burden in all these three mutation types differed between the 11
162 countries ($p\text{-value} < 2 \times 10^{-23}$, $p\text{-value} < 2 \times 10^{-14}$, $p\text{-value} < 6 \times 10^{-14}$, for SBSs, DBSs, and IDs,
163 respectively). In particular, the burden of all mutation types was elevated in Romania compared

164 to other countries (**Extended Data Fig. 1**). Principal Component Analysis (PCA) performed on
165 the proportions of the six primary SBS mutation classes (C>A, C>G, C>T, T>A, T>C, T>G) in
166 each sample identified a distinct cluster of mainly Romanian and Serbian cases and a further
167 cluster of mainly Japanese cases (**Extended Data Fig. 2**). The results, therefore, clearly
168 demonstrate geographical variation of somatic mutation loads and patterns in ccRCC.

169

170 To investigate the mutational processes contributing to the geographical variation in mutation
171 burdens we extracted mutational signatures and estimated the contribution of each signature
172 to each ccRCC genome. Ten signatures with strong similarity to a reference signature in the
173 Catalogue of Somatic Mutations in Cancer (COSMIC) database were extracted: SBS1, due to
174 deamination of 5-methylcytosine¹³; SBS2 and SBS13, due to cytosine deamination by
175 Apolipoprotein B mRNA-editing enzyme, catalytic polypeptide-like (APOBEC) DNA editing
176 enzymes¹³; SBS4, due to tobacco smoke mutagens¹⁴; SBS5, due to an endogenous
177 mutational process in which mutations accumulate with age¹⁴; SBS12, of unknown cause;
178 SBS18, due to DNA damage by reactive oxygen species¹⁴; SBS21 and SBS44, due to
179 defective DNA mismatch repair^{14,15}; and SBS22, due to Aristolochic acid exposure^{16,17}.

180

181 Five further SBS signatures were identified but could not be well described by the COSMIC
182 catalogue (**Fig. 2; Supplementary Table 5**). SBS40a, SBS40b and SBS40c were present in
183 most ccRCC accounting for, on average, ~30%, ~20%, and ~3% of mutations respectively
184 (**Fig. 2b**). Combined, they closely resemble the previously reported SBS40 (0.97 cosine
185 similarity), suggesting that the large number of ccRCC whole genomes analyzed here provides
186 the power to separate the constituent component signatures of SBS40. SBS40 was previously
187 reported frequently, and at high levels, in kidney cancer, but also in other cancers, and is of
188 unknown etiology. Like the composite SBS40, SBS40a is present in multiple cancer types.
189 However, SBS40b and SBS40c are largely restricted to ccRCC (**Supplementary Note 1**).
190 SBS_H was found in a single case and SBS_I is related to Aristolochic acid exposure (see

191 below; SBS_I has been renamed as SBS22b). Analysis of all other types of mutational
192 signatures is presented in Supplementary results.

193

194 The mutation burdens of multiple SBS mutational signatures varied between the 11 countries.
195 SBS22 is thought to be caused by Aristolochic acids, mutagenic derivatives of plants of the
196 Aristolochia genus which are carcinogenic and also cause Balkan endemic nephropathy
197 (BEN), a kidney disease prevalent in areas adjacent to the Danube in Southeastern Europe¹⁸.
198 SBS22 has previously been found in ccRCC, other urothelial tract cancers, and hepatocellular
199 carcinomas from Romania^{6,19} and various countries in East and South-East Asia^{16,17,20}. In this
200 study, SBS22 was present in high proportions of ccRCC from Romania (45/64, 70%), Serbia
201 (16/69, 23%), and Thailand (3/5, 60%), often with very high mutation burdens. The presence
202 of SBS22 was strongly correlated with that of new signatures SBS_I, DBS_D, and ID_C
203 (**Extended Data Fig. 3-5**) which are, therefore, also probably due to Aristolochic acid
204 exposure. SBS_I, like SBS22, is composed predominantly of T>A mutations. The signature
205 identified previously as SBS22, has therefore been renamed SBS22a, and SBS_I has been
206 named SBS22b. The two signatures may be due to different subsets of Aristolochic acids,
207 and/or to different metabolites, which induce slightly different mutational patterns. All these
208 signatures exhibited their highest mutation loads away from recognized BEN zones (**Fig. 3**)
209 indicating that a substantial proportion of the population over a wide geographical area of
210 Eastern Europe has been subject to mutagenesis due to Aristolochic acid exposure. The
211 sources of this exposure are uncertain.

212

213 SBS12 was present in 72% of Japanese and 2% of non-Japanese ccRCC (p -
214 value= 4.7×10^{-78}) (**Extended Data Fig. 6h**). Compared to the mutation burdens imposed by
215 Aristolochic acid in ccRCC, SBS12 contributed modest mutation loads. SBS12 is composed
216 predominantly of T>C substitutions and exhibits strong transcriptional strand bias with more
217 T>C mutations on the transcribed than untranscribed strands of protein coding genes.
218 Transcriptional strand bias is typically caused by activity of transcription-coupled nucleotide

219 excision repair acting on bulky DNA adducts due to exogenous mutagenic exposures such as
220 tobacco smoke chemicals¹⁴, ultraviolet light¹⁴, Aristolochic acids¹⁶, and aflatoxins²¹. Assuming
221 that transcription-coupled repair of DNA adducts is responsible for the SBS12 strand bias, the
222 adducts are likely on adenine. SBS12 was previously reported in hepatocellular carcinomas^{5,14}
223 and additional analysis of existing datasets revealed strong SBS12 enrichment in
224 hepatocellular carcinomas from Japan when compared to other countries ($p\text{-value}=3.8 \times 10^{-15}$;
225 **Supplementary Note**). The results, therefore, indicate that exposure to an agent contributing
226 SBS12 mutations to kidney and liver cancer is common in Japan and rare elsewhere. The
227 agent responsible for SBS12 is unknown. Although population-specific endogenous production
228 of the mutagen cannot be excluded, the precedents provided by other mutational signatures
229 with strong transcriptional strand bias suggest that it is likely of exogenous origin. A
230 polymorphism in aldehyde dehydrogenase 2 known to impair metabolism of alcohol to
231 aldehydes and common in Japan did not associate with levels of SBS12 (although power is
232 limited due to the relatively small number of Japanese cases).

233

234 SBS40a, SBS40b, and SBS40c were present in ccRCC from all 11 countries. The country-
235 specific average mutation burdens of SBS40a and SBS40b positively associated with country-
236 specific ASRs of kidney cancer incidence ($p\text{-value}=0.0022$ and $p\text{-value}=5.1 \times 10^{-18}$,
237 respectively; **Extended Data Fig. 7a; Fig. 4a**), with the highest mutation loads in the Czech
238 Republic and Lithuania. Indeed, when excluding the outlier effect of Romania and Serbia,
239 SBS40b was largely responsible for association of country-specific average total SBS burdens
240 with kidney cancer ASR ($p\text{-value}=6 \times 10^{-5}$; **Extended Data Fig. 8a**). Kidney cancer incidence
241 rates also vary between the regions of the Czech Republic and SBS40b mutation burdens
242 differed significantly between these ($p\text{-value}=0.011$; **Fig 4b,c**), with the highest attribution in
243 the highest risk region. SBS40b exhibits modest transcriptional strand bias and, assuming that
244 transcription-coupled repair of DNA adducts is responsible, the adducts underlying SBS40b
245 are likely on pyrimidines. Insertion and deletion (indel) signatures ID5 and ID8, which together
246 contributed ~60% of the indel mutation burden on average, were also strongly associated with

247 country-specific kidney cancer ASR (p -value= 1.3×10^{-10} and p -value= 6.2×10^{-9} , respectively,
248 **Extended Data Fig. 7b,c**). Signatures ID5 and ID8 correlated with each other (0.78), as well
249 as with SBS40b (0.81 and 0.74, respectively) indicating that they likely all constitute products
250 of the same underlying mutational process. Thus, the burdens of the full complement of
251 mutation types generated by this mutational process correlate with age-adjusted kidney cancer
252 incidence rates.

253

254 To investigate potential mutagenic agents underlying these geographically variable signatures,
255 an untargeted metabolomics screen of plasma was conducted on 901 individuals in the study,
256 from all countries except Japan (**Methods**). 2,392 metabolite features were obtained, including
257 944 independent peaks ($r < 0.85$). Three features were associated with SBS4 (**Supplementary**
258 **Table 13**), with two identified as hydroxycotinine (p -value= 2.9×10^{-9}) and cotinine (p -value= 1.9
259 $\times 10^{-5}$), two major metabolites of nicotine²². Eight features were associated with SBS40b
260 (**Supplementary Table 13**). One feature was identified as N,N,N-trimethyl-L-alanyl-L-proline
261 betaine (TMAP; p -value= 1.2×10^{-5}), increased levels of which correlate strongly with reduced
262 kidney function²³. Other established measures of kidney function, including cystatin C and
263 creatinine, were correlated with TMAP (p -value = 2.5×10^{-30} and 1.7×10^{-69} , respectively) and
264 also showed evidence of positive association with SBS40b (p -value=0.023 and 0.058,
265 respectively). Thus, exposure to the mutagenic agent responsible for SBS40b is associated
266 with reduced kidney function. No recognized metabolome features were significantly
267 associated with any other signatures.

268

269 A total of 1913 “driver” mutations were found in 136 genes including *VHL*, *PBRM1*, *SETD2*
270 and *BAP1*, the known frequently mutated cancer genes in ccRCC (**Methods**) (**Fig. 5a**)^{10,24}. The
271 frequencies of mutations in these genes were consistent across countries (**Fig. 5b**). The
272 spectrum of all driver mutations in ccRCC with Aristolochic acid exposure (**Methods**) was
273 enriched in T>A mutations compared to non-exposed cases (25% vs 13%, p -value=0.0062,
274 **Fig. 5c,d**) with similar enrichment specifically in *VHL* mutations (30% vs 16%; **Fig. 5e,f**), and

275 in the whole exome (27% in exposed compared to 12% in unexposed cases). Thus genome-
276 wide Aristolochic acid mutagenesis has contributed in a proportionate fashion to generation of
277 driver mutations in Aristolochic acid-exposed ccRCC. SBS12 did not show statistically
278 significant enrichment in drivers in exposed cases, possibly due to the smaller numbers of
279 SBS12 exposed cases and the lower mutation burden conferred. SBS40b also did not show
280 statistically significant enrichment probably due to the ubiquitous exposure and its relatively
281 flat and featureless mutation profile.

282

283 Exogenous mutagenic exposures that ultimately cause cancer may be present during the early
284 stages of evolution of cancer clones. To time mutagenic exposures, the contribution of each
285 mutational signature to mutations in the primary clone (relatively early) and to mutations in
286 subclones (relatively late) were estimated^{25,26} (**Methods**). All signatures of the putative
287 exogenous mutagenic exposures observed in ccRCC were present at relatively early stages
288 of cancer development, consistent with exposures to normal cells. SBS12, SBS22b, and
289 SBS40b showed higher activities in main clones compared to subclones (q-value=0.04, q-
290 value=0.02, q-value= 2.3×10^{-5} , respectively) (**Extended Data Fig. 9**) and SBS22a showed no
291 significant difference^{16,17}. By contrast, signatures due to endogenous mutational processes
292 including APOBEC DNA editing (SBS13) and oxidative damage (SBS18), were enriched in
293 subclones (q-value= 1.6×10^{-4} , q-value= 3.2×10^{-7} , respectively).

294

295 Established or suspected risk factors for ccRCC include age, tobacco smoking, obesity,
296 hypertension, diabetes, and environmental exposure to PFAS compounds³. Total SBS, DBS,
297 and ID mutation burdens associated with age, as did SBS1, SBS4, SBS5, SBS40a, SBS40b,
298 SBS22a, SBS22b, DBS2, ID1, ID5, and ID8. Total SBS (p-value= 2.8×10^{-5}), DBS (p-
299 value= 3.5×10^{-3}) and ID (p-value= 1.1×10^{-4}) mutation burdens also associated with sex, with
300 males having higher mutation burdens than females, and with SBS40b showing a similar
301 association (p-value= 7.6×10^{-5}). Associations with tobacco smoking were observed for SBS4
302 (p-value= 5.5×10^{-6}) and DBS2 (p-value= 2.3×10^{-7}), both known to be caused by tobacco

303 carcinogens^{27,28}, but associations of particular mutational signatures with other ccRCC risk
304 factors were not observed (**Supplementary Tables 8-9**). To complement this analysis of
305 observational data, associations between polygenic risk scores for known ccRCC risk factors
306 and mutational signatures^{29,30} were examined (**Methods**). Consistent with the observational
307 data, no associations were found between genetically inferred risk factors and mutational
308 signatures except for tobacco smoking and DBS2 (p-value=0.008; **Supplementary Table 10**).

309

310 **DISCUSSION**

311 Somatic mutations in the genomes of 962 ccRCC patients from 11 countries indicate the
312 existence of multiple, widespread mutational processes exhibiting substantial geographical
313 variation in their contributions to ccRCC mutation loads. The results contrast with those from
314 552 esophageal squamous carcinomas from eight countries with widely different esophageal
315 carcinoma incidence rates in which geographical differences in mutation burdens or signatures
316 were not observed³¹. Together the studies implicate both geographically variable mutagenic
317 and non-mutagenic carcinogenic exposures contributing to global cancer incidence. Indeed,
318 the presence of mutational signatures associated with tobacco smoking but absence of
319 signatures associated with other known ccRCC risk factors, such as obesity and hypertension,
320 suggests that the latter may be mediated by non-mutagenic processes and, therefore, that
321 both classes of carcinogen contribute to the development of ccRCC.

322

323 The existence, identity, and carcinogenic effect of some of the agents underlying these
324 mutational processes are known. Aristolochic acids are believed to cause SBS22a/b and its
325 associated signatures and this study suggests that the geographical extent and proportion of
326 the population acquiring mutations in Eastern Europe may be greater than previously
327 anticipated. The sources of the Aristolochic acid exposure, the manner by which they are
328 ingested and whether the exposure continues today are uncertain and further definition of the
329 source and extent of this exposure is required in order to provide a foundation for public health
330 action.

331

332 The existence of the mutagenic exposures underlying SBS12 and SBS40b were not previously
333 suspected, and their causative agents are unknown. Based on current information, the
334 exposure causing SBS12 is restricted to Japan. However, larger studies are now indicated to
335 explore the geographical extent of exposure in Japan and neighboring countries, and the
336 proportions of their populations developing mutations. In the first instance this will be
337 achievable by further sequencing of kidney and hepatocellular cancer genomes. However,
338 studies of normal tissues, using recently reported sequencing methods allowing detection of
339 somatic mutations in normal cells³², and particularly relatively accessible ones such as cells in
340 urine that can be prospectively collected, may enable large population-based studies providing
341 better characterization of the exposure and its consequences.

342

343 In contrast to Aristolochic acid and the agent causing SBS12, the exposure underlying SBS40b
344 appears to be globally ubiquitous. It predominantly causes mutations in ccRCC, with much
345 lower burdens in other cancer types, and generates mutation loads correlating strongly with
346 age and sex. There are few clues as to its origin or nature.

347

348 The incidence rates of ccRCC vary ~eightfold across the eleven countries from which ccRCCs
349 were sequenced. A strong positive correlation ($p\text{-value}=5.5 \times 10^{-18}$) was found between the
350 average mutation loads attributable to SBS40b in each country (and also those of ID5 and ID8
351 which are correlated with SBS40b) and incidence of kidney cancer within each country. This
352 correlation reflects approximately a tripling of average country-specific SBS40b mutation loads
353 (a difference of ~1000 mutations) in parallel with the eightfold increase of country-specific ASR.

354

355 SBS40b mutation burdens also positively correlated with biomarkers of impaired kidney
356 function, reminiscent of the nephrotoxic effects of Aristolochic acids in Balkan endemic
357 nephropathy. It is possible that the increased SBS40b somatic mutation load itself engenders
358 this reduction in renal function. However, studies of other normal tissues suggest that they are

359 generally tolerant of elevated mutation burdens, except for manifesting a higher incidence of
360 neoplasia^{33,34}. It is also possible that the agent underlying SBS40b is directly nephrotoxic, for
361 example by engendering DNA damage and a response to it, and that the mutations it
362 generates are immaterial to kidney function. It is also conceivable, however, that impaired renal
363 function, potentially due to many different causes, results in a metabolic state which itself
364 causes the elevated SBS40b mutation load. Whatever the mutational process underlying
365 SBS40b, it is plausible that it contributes to the geographical variation in the age standardized
366 rates of kidney cancer incidence rates. It is of public health interest to determine the cause of
367 SBS40b and, hence, to consider whether the exposure can be mitigated, potentially with
368 concomitant reduction in global ccRCC incidence rates.

369

370 Finally, it is notable that overall tumor mutation burden did not vary substantially between
371 countries, with most variation being due to mutations linked to Aristolochic acid, and had only
372 a weak association with overall kidney cancer incidence (**Extended Data Fig. 1 and 8**). Along
373 with an absence of any association between several known risk factors for ccRCC and
374 mutation burden, in particular for obesity and hypertension, these results provide further
375 evidence for a model of cancer development where mutations are essential but additional
376 factors affect the expansion of a mutated clone and thus the chance of it progressing into
377 cancer³⁵. Further efforts at defining how lifestyle and environmental exposures contribute to
378 cancer development will therefore require a greater understanding of both the causes of the
379 mutations in cell clones in normal tissue, and the further promotion of such mutant clones by
380 non-mutagenic processes.

381

382 REFERENCES

- 383 1. Brennan, P. & Davey-Smith, G. Identifying Novel Causes of Cancers to Enhance
384 Cancer Prevention: New Strategies Are Needed. *JNCI: Journal of the National Cancer*
385 *Institute* **114**, 353–360 (2022).
- 386 2. Hsieh, J. J. *et al.* Renal cell carcinoma. *Nat Rev Dis Primers* **3**, 17009 (2017).

- 387 3. Shearer, J. J. *et al.* Serum Concentrations of Per- and Polyfluoroalkyl Substances and
388 Risk of Renal Cell Carcinoma. *J Natl Cancer Inst* **113**, 580–587 (2021).
- 389 4. Koh, G., Degasperi, A., Zou, X., Momen, S. & Nik-Zainal, S. Mutational signatures:
390 emerging concepts, caveats and clinical applications. *Nat Rev Cancer* **21**, 619–637
391 (2021).
- 392 5. Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer.
393 *Nature* **578**, 94–101 (2020).
- 394 6. Scelo, G. *et al.* Variation in genomic landscape of clear cell renal cell carcinoma
395 across Europe. *Nat Commun* **5**, 5135 (2014).
- 396 7. Mitchell, T. J. *et al.* Timing the Landmark Events in the Evolution of Clear Cell Renal
397 Cell Cancer: TRACERx Renal. *Cell* **173**, 611-623.e17 (2018).
- 398 8. Campbell, P. J. *et al.* Pan-cancer analysis of whole genomes. *Nature* **578**, 82–93
399 (2020).
- 400 9. Degasperi, A. *et al.* A practical framework and online tool for mutational signature
401 analyses show intertissue variation and driver dependencies. *Nat Cancer* **1**, 249–263
402 (2020).
- 403 10. The Cancer Genome Atlas Research Network. Comprehensive molecular
404 characterization of clear cell renal cell carcinoma. *Nature* **499**, 43–49 (2013).
- 405 11. Mutographs Cancer Grand Challenge. <https://cancergrandchallenges.org/teams>.
- 406 12. Sung, H. *et al.* Global Cancer Statistics 2020: GLOBOCAN Estimates of Incidence
407 and Mortality Worldwide for 36 Cancers in 185 Countries. *CA Cancer J Clin* **71**, 209–
408 249 (2021).
- 409 13. Nik-Zainal, S. *et al.* Mutational Processes Molding the Genomes of 21 Breast
410 Cancers. *Cell* **149**, 979–993 (2012).
- 411 14. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature*
412 **500**, 415–421 (2013).
- 413 15. Drost, J. *et al.* Use of CRISPR-modified human stem cell organoids to study the origin
414 of mutational signatures in cancer. *Science (1979)* **358**, 234–238 (2017).

- 415 16. Hoang, M. L. *et al.* Mutational Signature of Aristolochic Acid Exposure as Revealed by
416 Whole-Exome Sequencing. *Sci Transl Med* **5**, (2013).
- 417 17. Poon, S. L. *et al.* Genome-wide mutational signatures of aristolochic acid and its
418 application as a screening tool. *Sci Transl Med* **5**, 197ra101 (2013).
- 419 18. Grollman, A. P. Aristolochic acid nephropathy: Harbinger of a global iatrogenic
420 disease. *Environ Mol Mutagen* **54**, 1–7 (2013).
- 421 19. Turesky, R. J. *et al.* Aristolochic acid exposure in Romania and implications for renal
422 cell carcinoma. *Br J Cancer* **114**, 76–80 (2016).
- 423 20. Wang, X.-M. *et al.* Integrative genomic study of Chinese clear cell renal cell carcinoma
424 reveals features associated with thrombus. *Nat Commun* **11**, 739 (2020).
- 425 21. Huang, M. N. *et al.* Genome-scale mutational signatures of aflatoxin in cells, mice, and
426 human tumors. *Genome Res* **27**, 1475–1486 (2017).
- 427 22. Dempsey, D. *et al.* Nicotine metabolite ratio as an index of cytochrome P450 2A6
428 metabolic activity. *Clin Pharmacol Ther* **76**, 64–72 (2004).
- 429 23. Velenosi, T. J. *et al.* Untargeted metabolomics reveals N, N, N-trimethyl-L-alanyl-L-
430 proline betaine (TMAP) as a novel biomarker of kidney function. *Sci Rep* **9**, 6831
431 (2019).
- 432 24. Sato, Y. *et al.* Integrated molecular analysis of clear-cell renal cell carcinoma. *Nat*
433 *Genet* **45**, 860–867 (2013).
- 434 25. Nik-Zainal, S. *et al.* The life history of 21 breast cancers. *Cell* **149**, 994–1007 (2012).
- 435 26. Dentre, S. C., Wedge, D. C. & van Loo, P. Principles of Reconstructing the Subclonal
436 Architecture of Cancers. *Cold Spring Harb Perspect Med* **7**, (2017).
- 437 27. Nik-Zainal, S. *et al.* The genome as a record of environmental exposure. *Mutagenesis*
438 *gev073* (2015) doi:10.1093/mutage/gev073.
- 439 28. Kucab, J. E. *et al.* A Compendium of Mutational Signatures of Environmental Agents.
440 *Cell* **177**, 821-836.e16 (2019).

- 441 29. Gabriel, A. A. G. *et al.* Genetic Analysis of Lung Cancer and the Germline Impact on
442 Somatic Mutation Burden. *JNCI: Journal of the National Cancer Institute* **114**, 1159–
443 1166 (2022).
- 444 30. Liu, Y., Gusev, A., Heng, Y. J., Alexandrov, L. B. & Kraft, P. Somatic mutational
445 profiles and germline polygenic risk scores in human cancer. *Genome Med* **14**, 14
446 (2022).
- 447 31. Moody, S. *et al.* Mutational signatures in esophageal squamous cell carcinoma from
448 eight countries with varying incidence. *Nat Genet* **53**, 1553–1563 (2021).
- 449 32. Abascal, F. *et al.* Somatic mutation landscapes at single-molecule resolution. *Nature*
450 **593**, 405–410 (2021).
- 451 33. Martincorena, I. *et al.* High burden and pervasive positive selection of somatic
452 mutations in normal human skin. *Science (1979)* **348**, 880–886 (2015).
- 453 34. Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus with age.
454 *Science (1979)* **362**, 911–917 (2018).
- 455 35. Fowler, J. C. & Jones, P. H. Somatic Mutation: What Shapes the Mutational
456 Landscape of Normal Epithelia? *Cancer Discov* **12**, 1642–1655 (2022).

457

458 **FIGURE AND TABLE LEGENDS**

459 **Fig. 1: Eleven participating countries and estimated age-standardized incidence rates**
460 **of clear cell renal cell carcinomas.**

461 Incidence of clear cell renal cell carcinomas (ccRCC), men and women combined, age-
462 standardized incidence rates (ASR) per 100,000, data from GLOBOCAN 2020. Markers
463 indicate countries included in this study (number of participating ccRCC patients per
464 country).

465

466 **Table 1. Summary of clear cell renal cell carcinomas risk factors included in this**
467 **study.**

468

469 **Fig. 2: Single base substitution signature operative in clear cell renal cell carcinomas.**

470 **(a)** TMB plot showing the frequency and mutations per Mb for each of the decomposed SBS
471 signatures. **(b)** Average relative attribution for single base substitution (SBS) signatures
472 across countries. Signatures contributing less than 5% on average are grouped in the
473 'Others' category, apart from SBS12 and AA-related signatures SBS22a and SBS22b. '<95%
474 confidence' category accounts for the proportion of mutation burden which could not be
475 assigned to any signature with confidence level of at least 95%. **(c)** Decomposed signatures,
476 including reference COSMIC signatures as well as *de novo* signatures not decomposed into
477 COSMIC reference signatures.

478

479 **Fig. 3: Geospatial analysis of Aristolochic acid-related SBS signatures.**

480 Distribution of Romanian and Serbian cases with known residential history, along with the
481 summed levels of SBS22a and SBS22b attributions (per-case and regional estimate), with
482 respect to the Balkan endemic nephropathy (BEN) areas. White circles represented cases
483 with no detected activity of SBS22a and SBS22b.

484

485 **Fig. 4: Association of SBS40b signature attribution with incidence of kidney cancer.**

486 **(a)** Number of mutations attributed to signature SBS40b against age-standardized incidence
487 rates (ASR) of kidney cancer in each of the eleven countries represented in the cohort. **(b)**
488 Number of mutations attributed to signature SBS40b in four regions of Czech Republic
489 against ASR of kidney cancer in each region. **(c)** Levels of attribution of SBS40b signature
490 within Czech Republic, with bar plots showing the number of cases for each quartile of
491 SBS40b attribution across Prague, Olomouc, Ceske Budejovice, and Brno regions.

492

493 **Fig. 5: Driver mutation analysis in clear cell renal cell carcinomas.**

494 **(a)** Frequency of driver genes in the cohort. Only genes mutated in at least 10 cases are
495 shown. **(b)** Frequency of driver genes across countries. **(c)** SBS-96 mutational spectra of all
496 driver mutations in ccRCC for Aristolochic acid (AA)-exposed and unexposed cases. **(d)**

- 497 Percentage of T>A driver mutations in AA-exposed and unexposed cases. **(e)** SBS-96
- 498 mutational spectra of VHL mutations in ccRCC for AA-exposed and unexposed cases. **(f)**
- 499 Percentage of T>A VHL mutations in AA-exposed and unexposed cases.

500 **ONLINE METHODS**

501 **Recruitment of cases and informed consent**

502 The International Agency for Research on Cancer (IARC/WHO) coordinated case recruitment
503 through an international network of over 40 collaborators from the 11 participating countries
504 (**Table1; Supplementary table 12**). The inclusion criteria for patients were ≥ 18 years of age
505 (ranging from 23 to 87, with a mean of 60 and a standard deviation of 12), confirmed diagnosis
506 of primary ccRCC and no prior cancer treatment. Informed consent was obtained for all
507 participants. Patients were excluded if they had any condition that could interfere with their
508 ability to provide informed consent or if there were no means of obtaining adequate tissues or
509 associated data as per the protocol requirements. Ethical approvals were first obtained from
510 each Local Research Ethics Committee and Federal Ethics Committee when applicable, as
511 well as from the IARC Ethics Committee.

512

513 **Bio-samples, data collection, and expert pathology review**

514 Dedicated standard operating procedures, following guidelines from the International Cancer
515 Genome Consortium (ICGC), were designed by IARC/WHO to select appropriate case series
516 with complete biological samples and exposure information as described previously¹
517 (**Supplementary Table 12**). In brief, for all case series included, anthropometric measures
518 were taken, together with relevant information regarding medical and familial history.
519 Comparable smoking and alcohol history was available from all centers. Detailed
520 epidemiological information on residential history was collected in Czech Republic, Romania,
521 and Serbia. Potential limitations of using retrospective clinical data collected using different
522 protocols from different populations were addressed by a central data harmonization to ensure
523 a comparable group of exposure variables (**Supplementary Table 12**). All patient related data
524 as well as clinical, demographical, lifestyle, pathological and outcome data were
525 pseudonymized locally through the use a dedicated alpha-numerical identifier system before
526 being transferred to IARC/WHO central database.

527 Original diagnostic pathology departments provided diagnostic histological details of
528 contributing cases through standard abstract forms. IARC/WHO centralized the entire
529 pathology workflow and coordinated a centralized digital pathology examination of the frozen
530 tumor tissues collected for the study as well as formalin-fixed, paraffin-embedded (FFPE)
531 sections when available, via a web-based report approach and dedicated expert panel
532 following standardized procedures as described previously¹. A minimum of 50% viable tumor
533 cells was required for eligibility to whole genome sequencing.

534 In summary, frozen tumor tissues were first examined to confirm the morphological type and
535 the percentage of viable tumor cells. A random selection of tumor tissues was independently
536 evaluated by a second pathologist. Enrichment of tumor component was performed by
537 dissection of non-tumoral part, if necessary.

538

539 **DNA extraction**

540 Extraction of DNA from fresh frozen tumor and matched blood samples was centrally
541 conducted at IARC/WHO except for Japan, which performed DNA extractions at the local
542 center following a similarly standardized DNA extraction procedure. Of the cases which
543 proceeded to the final analysis ($n=962$), germline DNA was extracted from either buffy-coat,
544 whole blood, or from adjacent normal tissue (*viz.*, samples from Japan) using previously
545 described protocols and methods¹.

546

547 **Whole genome sequencing**

548 In total, 1583 renal cell carcinoma cases were evaluated, with 1267 confirmed as ccRCC
549 cases. 116 (9%) were excluded due to insufficient viable tumor cells (pathology level), or
550 inadequate DNA (tumor or germline). DNA from 1151 cases was received at the Wellcome
551 Sanger Institute for whole genome sequencing. Fluidigm SNP genotyping with a custom panel
552 was performed to ensure that each pair of tumor and matched normal samples originated from
553 the same individual. Whole genome sequencing (150bp paired end) was performed on the
554 Illumina NovaSeq 6000 platform with target coverage of 40X for tumors and 20X for matched

555 normal tissues. All sequencing reads were aligned to the GRCh38 human reference genome
556 using Burrows-Wheeler-MEM (v0.7.16a and v0.7.17). Post-sequencing QC metrics were
557 applied for total coverage, evenness of coverage and contamination. Cases were excluded if
558 coverage was below 30X for tumor or 15X for normal tissue. For evenness of coverage, the
559 median over mean coverage (MoM) score was calculated. Tumors with MoM scores outside
560 the range of values determined by previous studies to be appropriate for whole genome
561 sequencing (0.92 – 1.09) were excluded. Conpair² (<https://github.com/nygenome/Conpair>)
562 was used to detect contamination, cases were excluded if the result was greater than 3%³. A
563 total of 962 cases passed all criteria and were included in subsequent analysis.

564

565 **Somatic variant calling**

566 Variant calling was performed using the standard Sanger bioinformatics analysis pipeline
567 (<https://github.com/cancerit>). Copy number profiles were determined first using the algorithms
568 ASCAT⁴ and BATTENBERG⁵, where tumor purity allowed. SNV were called with
569 cgpCaVEMan⁶, indels were called with cgpPINDEL⁷, and structural rearrangements were
570 called using BRASS. CaVEMan and BRASS were run using the copy number profile and purity
571 values determined from ASCAT where possible (complete pipeline, n=857). Where tumor
572 purity was insufficient to determine an accurate copy number profile (partial pipeline, n=105),
573 CaVEMan and BRASS were run using copy number defaults and an estimate of purity
574 obtained from ASCAT/BATTENBERG. For SNV additional filters (ASRD \geq 140 and CLPM
575 $=0$) were applied to remove potential false positive calls. To further exclude the possibility of
576 caller specific artefacts being included in the analysis, a second variant caller, Strelka2, was
577 run for SNVs and indels^{1,8}. Only variants called by both the Sanger variant calling pipeline and
578 Strelka2 were included in subsequent analysis.

579

580 **Validation of Japanese sequencing**

581 The matched normal tissue sequenced was blood for all countries with the exception of Japan,
582 where adjacent normal kidney was used. As Japan displayed an enrichment of SBS12,

583 matched blood was obtained from 28 of the 36 patients to confirm that the source of the
584 matched normal tissue was not influencing the result. In all cases, the mutational spectra of
585 Japanese ccRCC generated using either blood or adjacent normal kidney matched each other
586 with a cosine similarity of >0.99.

587

588 **Generation of mutational matrices**

589 Mutational matrices for single base substitutions (SBS), doublet base substitutions (DBS) and
590 small insertions and deletions (ID) were generated using SigProfilerMatrixGenerator
591 (<https://github.com/AlexandrovLab/SigProfilerMatrixGenerator>) with default options (v1.2.12)⁹.

592

593 **Mutational signature analysis**

594 Mutational signatures were extracted using two algorithms, SigProfilerExtractor
595 (<https://github.com/AlexandrovLab/SigProfilerExtractor>), based on nonnegative matrix
596 factorization, and mSigHdp¹⁰ (<https://github.com/steverozen/mSigHdp>), based on the
597 Bayesian hierarchical Dirichlet process. For SigProfilerExtractor, *de novo* mutational
598 signatures were extracted from each mutational matrix using SigProfilerExtractor with
599 nndsvd_min initialization (NMF_init="nndsvd_min") and default parameters (v1.1.9)¹¹. Briefly,
600 SigProfilerExtractor deciphers mutational signatures by first performing Poisson resampling of
601 the original matrix with additional renormalization (based on a generalized mixture model
602 approach) of hypermutators to reduce their effect on the overall factorization¹¹. Nonnegative
603 matrix factorization (NMF) was performed using initialization with nonnegative singular value
604 decomposition and by applying the multiplicative update algorithm using the Kullback–Leibler
605 divergence as an objective function¹¹. NMF was applied with factorizations between $k=1$ and
606 $k=20$ signatures; each factorization was repeated 500 times¹¹. *De novo* single base
607 substitution mutational signatures were extracted with SigProfilerExtractor for both SBS-288
608 and SBS-1536 contexts⁹. The results were largely concordant with the SBS-1536 *de novo*
609 signatures allowing additional separation of mutational processes, therefore the SBS-1536 *de*
610 *nov*o signatures were taken forward for further analysis (**Supplementary Table 2**). Mutational

611 signatures for DBS and ID were extracted in DBS-78 and ID-83 contexts respectively
612 (**Supplementary Tables 3 & 4**). Where possible, SigProfilerExtractor matched each *de novo*
613 extracted mutational signature to a set of previously identified COSMIC signatures¹², for SBS-
614 1536 signatures this requires collapsing the 1536 classification into the standard 96
615 substitution type classification with six mutation classes having single 3' and 5' sequence
616 contexts (Supplementary Table 5). This step makes it possible to distinguish between *de novo*
617 signatures which can be explained by a combination of the known catalog of mutational
618 process (which have not been completely separated during the extraction), and those which
619 have not been previously identified. mSigHdp extraction of SBS-96 and ID-83 signatures was
620 performed using the suggested parameters and using the country of origin to construct the
621 hierarchy. SigProfilerExtractor's decomposition module was subsequently used to match
622 mSigHdp *de novo* signatures to previously identified COSMIC signatures¹². Further details on
623 the comparison of results between SigProfilerExtractor and mSigHdp and decomposition of *de*
624 *nov*o signatures into COSMIC reference signatures can be found in the supplementary note.

625

626 **Attribution of activities of mutational signatures**

627 The *de novo* (SigProfiler) and COSMIC signature (SigProfiler and mSigHdp) activities were
628 attributed for each sample using the MSA signature attribution tool (v2.0,
629 <https://gitlab.com/s.senkin/MSA>)¹³. For COSMIC attributions, only COSMIC reference
630 signatures, which were identified in the decomposition of *de novo* signatures, were included
631 in the panel for attribution, in addition to *de novo* signatures which could not be decomposed
632 into COSMIC reference. At its core, the tool utilizes the nonnegative least squares (NNLS)
633 approach minimizing the L2 distance between the input sample and the one reconstructed
634 using available signatures. To limit false positive attributions, automated optimization
635 procedure was applied by repeated removal of all signatures that do not increase the L2
636 similarity of a sample by >0.008 for SBS, >0.014 for DBS, and >0.03 for ID mutation types, as
637 suggested by simulations. These optimal penalties were derived using an optional parameter
638 (params.no_Cl_for_penalties = false) utilizing a conservative approach in calculation of

639 penalties. Finally, a parametric bootstrap approach was applied to extract 95% confidence
640 intervals for each attributed mutational signature activity.

641

642 **Driver mutations**

643 A dNdS approach was used to identify genes under positive selection in ccRCC¹⁴. The analysis
644 was performed both for the whole genome (q-value<0.01), and with restricted hypothesis
645 testing (RHT) for a panel of 369 known cancer genes¹⁴. Variants in any gene identified as
646 under positive selection in global dNdS or in the 369-cancer gene panel were assessed as
647 potential drivers¹⁴. Candidate driver mutations were annotated with the mode of action using
648 the Cancer Gene Census (<https://cancer.sanger.ac.uk/census>) and the Cancer Genome
649 Interpreter tool (<https://www.cancergenomeinterpreter.org>). Missense mutations were
650 assessed using the MutationMapper tool (http://www.cbiportal.org/mutation_mapper).
651 Variants were considered likely drivers if they met any of the following criteria: (i) Truncating
652 mutations in genes annotated as tumor suppressors; (ii) mutations annotated as likely or
653 known oncogenic in MutationMapper; (iii) truncating variants in genes with selection (q-
654 value<0.05) for truncating mutations assumed to be tumor suppressors and thus likely drivers;
655 (iv) missense variants in all genes under positive selection and with dN/dS ratios for missense
656 mutations above 5 (assuming 4 of every 5 missense mutations are drivers) labelled as likely
657 drivers; or (v) in-frame indels in genes under significant positive selection for in-frame indels.

658

659 **Evolutionary analysis**

660 Subclonal architecture reconstruction was performed using the DPCLust R package v2.2.8^{5,15},
661 after obtaining cancer cell fraction (CCF) estimates by dpclust3p v1.0.8
662 (<https://github.com/Wedge-lab/dpclust3p>) based on the variant allele frequency provided by
663 the somatic variant callers and the copy number profiles determined by the BATTENBERG
664 algorithm. Only tumors with at least 40% purity according to BATTENBERG were considered
665 for further evolutionary analysis. For each tumor with at least one subclone, the respective
666 somatic mutations were split into clonal and subclonal mutations using the most probable

667 cluster assignment for each mutation as per the DPCLust output. Mutations not assigned to a
668 cluster by DPCLust were removed from further analysis. Clusters centered at a CCF>1.5 and
669 ones where chromosome X contributed the highest number of mutations were deemed
670 artifactual, and the respective mutations were removed. Samples with a total number of clonal
671 or subclonal mutations below 256 were also removed. Additionally, samples with poor
672 separation between the clonal and subclonal distributions (e.g., subclone centered at a
673 CCF>0.80) were removed. Finally, only samples that had both a clone and at least one
674 subclone post-filtering were retained for further analysis. This yielded a total of 223 samples,
675 each with clonal and subclonal mutations. SigProfilerAssignment (v0.0.13)
676 (<https://github.com/AlexandrovLab/SigProfilerAssignment>) was used to identify the activity of
677 each mutational signature in each clone/subclone, and these activities were then normalized
678 by the total number of mutations belonging to the clone/subclone (i.e., clonal mutations were
679 not included in the subclone). A two-sided Wilcoxon Signed-Rank Test¹⁶ was used to assess
680 the differences in the relative activity of each mutational signature between the clones and
681 their respective subclones. P-values were corrected using the Benjamini-Hochberg
682 procedure¹⁷ and reported as q-values in the manuscript.

683

684 **Regressions**

685 Signature attributions were dichotomized into presence and absence using confidence
686 intervals, with presence defined as both lower and upper limits being positive, and absence as
687 the lower limit being zero. If a signature was present in at least 75% of cases (SBS1, SBS40a,
688 SBS40b, ID1, and ID5), it was dichotomized into above and below the median of attributed
689 mutation counts. The binary attributions served as dependent variables in logistic regressions,
690 and relevant risk factors were used as factorized independent variables. To adjust for
691 confounding factors, sex, age of diagnosis, country, and tobacco status were added as
692 covariates in regressions. The Bonferroni method was used to test for significant p-values (i.e.,
693 a total of 224 comparisons for regressions with signatures, and a total of 24 comparisons for
694 regressions with mutation burden). P-values reported are raw (not corrected). Regressions

695 with incidence of renal cancer were performed as linear regressions with mutation burdens or
696 signature attributions with confidence intervals not consistent with zero as a dependent
697 variable, and age-standardized rates (ASR) of renal cancer obtained from Global Cancer
698 Observatory (GLOBOCAN)¹⁸, sex and age of diagnosis as independent variables. ASR of renal
699 cancer for regions of Czech Republic were obtained from SVOD web portal¹⁹.

700

701 **Polygenic risk score (PRS) analysis of lifestyle risk factors**

702 In this analysis, we used the genome-wide association studies (GWAS) summary statistics
703 estimated in European populations for well-established risk factors for ccRCC. For tobacco
704 smoking status, we used results from the GSCAN consortium meta-analysis of smoking
705 initiation (ever vs never status)²⁰. For body mass index (BMI), the results of UK biobank (UKBB)
706 meta-analysis of continuous BMI were used²¹. GWAS summary statistics related to
707 hypertension, namely systolic blood pressure and diastolic blood pressure, as well as the ones
708 related to diabetes²², such as fasting glucose and fasting insulin were also obtained using
709 UKBB studies²³.

710

711 Since all the GWAS summary statistics used in the current work were based on European
712 populations, we used ADMIXTURE tool (v1.3.0)²⁴ and principal component analysis (PCA) to
713 infer the unsupervised cluster of individuals with European genetic background within ccRCC
714 cases. Hapmap SNPs (n=1,176,821 variants) were extracted from the ccRCC whole-genome
715 sequence genotype data. After basic quality control using PLINK (v1.9b, [www.cog-
716 genomics.org/plink/1.9/](http://www.cog-genomics.org/plink/1.9/)), 333 variants were removed due to missing genotype rate > 5%, 1,236
717 variants failed Hardy-Weinberg equilibrium test (p-values<10⁻⁸), and 18,702 variants had
718 MAF<1% in our cohort. Additionally, 3 ambiguous variants and 21,358 variants within regions
719 of long-range, high linkage disequilibrium (LD) in the human genome (hg38) were excluded.
720 After pruning for linkage disequilibrium, 143,727 variants remained in ccRCC genotype data.
721 The 1000 genome reference population genotype data (phase 3) for Europeans (N=489),
722 Africans (YRI, N=108) and East Asians (N=103 from China and 104 from Japan)

723 (<https://www.internationalgenome.org/data/>) were filtered and merged with ccRCC genotype
724 data based on the pruned set of variants present in both datasets. ADMIXTURE was run on
725 the merged genotype data with $k=3$, which would correspond to the three ancestral continental
726 population groups that likely reflect the participants of our study. The ccRCC cases with
727 European genetic fraction greater than 80% by the ADMIXTURE analysis were selected for
728 the polygenic risk scores (PRS) analyses. To complement the ADMIXTURE analysis, PCA
729 was run on the same samples.

730

731 The initial genotype data based on whole-genome sequence from 849 ccRCC cases with
732 European genetic background consisted of biallelic SNPs with MAF $>0.01\%$ (to exclude ultra-
733 rare variants; $N \sim 30$ million variants). After basic quality control, variants with missing
734 genotype rate of greater than 5% ($N=7,519,196$ variants) with strong deviation from Hardy-
735 Weinberg equilibrium ($p\text{-values} < 10^{-8}$, $N=220,862$) were excluded. For each GWAS trait, we
736 restricted our analyses to the biallelic SNPs with minor allele frequency (MAF) greater than 1%
737 in the 1000 genomes reference for European populations. For the selection of the independent
738 genome-wide significant hits ($p\text{-values} < 5 \times 10^{-8}$) of each GWAS summary statistic used to
739 generate the PRS, SNPs were clumped ($r^2=0.1$ within a LD window of 10 MB) using PLINK
740 (v1.9b, www.cog-genomics.org/plink/1.9/) based on the 1000 genomes European reference
741 population genotype data ($N=489$; ~ 10 million variants). Where a selected GWAS hit was not
742 found in ccRCC genotype data, we extracted proxies ($r^2 > 0.8$ in 1000 genomes) also present
743 in ccRCC dataset where possible (**Supplementary Table 11**). The variance of each genetic
744 trait explained by the genetic variants were calculated as previously suggested²⁵. PRS was
745 subsequently calculated as the sum of the individual's beta-weighted genotypes using PRSice-
746 2 software²⁶. Associations were estimated per standard deviation increase in the PRS, which
747 was normalized to have a mean of zero across ccRCC cases of European genetic ancestry.

748

749

750

751 **Untargeted metabolomics association with signatures**

752 Of the 962 subjects from the main analysis, 901 subjects were included in this sub-study – all
753 Japanese samples ($n=36$) as well as few cases from Czech Republic ($n=13$), Romania ($n=5$)
754 and Russia ($n=3$) were not included due to lack of available plasma samples. Samples were
755 randomized and analyzed as two independent analytical batches. Analysis was performed with
756 a UHPLC-QTOF-MS system that consisted of a 1,290 Binary LC and a 6,550 QTOF mass
757 spectrometer equipped with Jet Stream electrospray ionization source (Agilent Technologies),
758 using previously described methods²⁷. Pre-processing was performed using Profinder
759 10.0.2.162 and Mass Profiler Professional B.14.9.1 software (Agilent Technologies). A “Batch
760 recursive feature extraction (small molecules)” process was employed for samples and blanks
761 to find $[M+H]^+$ ions. The two batches were processed separately and the resulting features
762 were aligned in Mass Profiler Professional. Chromatographic peak areas were used as a
763 measurement of intensity.

764

765 A total of 2,392 features were detectable in at least one of the 901 samples. Features present
766 in only one of the two batches were filtered out. Recursive filtering elimination was applied to
767 decrease redundancy from highly correlated variables ($r \geq 0.85$, Pearson's r calculated before
768 any transformation/imputation) by selecting the features with least missing data within clusters
769 of features. A total of 944 features were included in the statistical analysis. Features were pre-
770 processed: missing values were replaced with 1/5 of the minimal value of the feature before
771 applying mean centering and Pareto scaling. Each feature was regressed against both *de novo*
772 and COSMIC signatures, adjusting for sex and age of diagnosis, as well as body mass index
773 (BMI) and technical factors (batch, acquisition order) that could impact chromatographic peak
774 area. Models for SBS22a and SBS22b were restricted to Romanian and Serbian samples to
775 find potential pathways of Aristolochic acid exposure in the Balkan region. Logistic models
776 were used for zero-inflated signatures ($\geq 30\%$ zeros) while quasi-Poisson regressions were
777 used for the least zero-inflated signatures (SBS1, SBS40a, and SBS40b). To derive specific
778 false detection rates, random variables were created from permutations of the initial features

779 and regressed against signatures in the same fashion as true features. Maximum p-value
780 thresholds from regressions with random features were compared to adjusted p-value
781 thresholds according to Bonferroni's procedure. The more conservative approach was used in
782 selecting features of interest. Random forest models were also used as cross-checking
783 multivariate models to assess the relative importance of each feature in explaining the
784 signature attribution. As with univariate models, regression models were used for the least
785 zero-inflated signatures (<30% of zeros) while classification models were used for all other
786 signatures, with restriction to Romanian and Serbian samples for SBS22a and SBS22b.
787 Importance was estimated from the total decrease in node impurities from splitting on the
788 variable, averaged over all trees. Node impurity was measured by the Gini index for
789 classification, and by residual sum of squares for regression. The significance of importance
790 metrics for Random Forest models were estimated by permuting the response variable
791 (<https://github.com/EricArcher/rfPermute>).

792

793 Features considered for identification, along with their highly correlated counterparts, were
794 searched in Human Metabolome Database (HMDB), LipidMaps, Metlin, and Kegg. Compound
795 identity was confirmed by comparison of retention times and MS/MS fragmentation against
796 chemical standards when available, or otherwise against reference MS/MS spectra. Since the
797 feature 240.1468@0.8929933 was strongly correlated with several features identified as
798 TMAP (Supplementary Table 13), the integration of these features was inspected and
799 corrected manually, and regressed against SBS40b using the same model applied to features
800 selected for analysis. Creatinine was identified among the features by matching its retention
801 time and MS/MS spectra against a reference standard and also regressed against SBS40b in
802 the same fashion as other metabolites. Estimation of correlation between metabolic features
803 was done using linear regression adjusting for batch and acquisition order.

804

805

806

807 **Targeted metabolomics analyses**

808 Circulating levels of PFAS (Per- and Polyfluorinated Substances) and cystatin C compounds
809 were investigated using targeted mass spectrometry-based methods as described
810 previously^{28,29}.

811

812 Out of the 962 subjects from the main analysis, plasma samples from 909 subjects (from all
813 countries except Japan) were randomized and sent frozen in dry ice to each respective
814 laboratory for analyses. Measurement of cystatin C from 906 subjects included its native form
815 and isoforms (3Pro-OH cystatin C, cystatin C-desS, 3Pro-OH cystatin C-desS and cystatin C-
816 desSSP) that were modeled individually and for the total concentration of cystatin C isoforms.
817 Measurements of PFAS compounds included PFOA (Perfluorooctanoic Acid; total, branch,
818 linear), PFOS (Perfluorooctanoic Acid; total, branch, linear), PFHxS (Perfluorohexane
819 sulfonate), PFNA (Perfluorononanoic acid), PFDA (Perfluorodecanoic acid), MePFOSAA (n-
820 methylperfluoro-1 octanesulfonamido acetic acid), EtPFOSAA (2-(N-Ethyl-perfluorooctane.
821 sulfonamido) acetic acid).

822

823 Multivariable quasi-Poisson (for the least sparse signatures SBS1, SBS40a and SBS40b) and
824 logistic regression were used to estimate the association between plasma concentrations of
825 the aforementioned substances and mutational signatures. All compounds were modeled
826 continuously (log₂-transformed) and categorically, with adjustments made by sex, age, date
827 of recruitment, country, BMI, tobacco and alcohol status in the case of PFAS molecules and
828 by sex, age and BMI, in the case of cystatin C.

829

830 **Geospatial analyses**

831 Geospatial analyses were performed to estimate the regional effect for signature attribution,
832 particularly for signatures thought to be from exogenous exposure (SBS40b – unknown – and
833 SBS22a/SB22b - Aristolochic acid). Residential history information was available for a large
834 proportion of cases from the countries of interest: Czech Republic for SBS40b and

835 Romania/Serbia for SBS22a/SBS22b. The 259 cases from Czech Republic within this study
836 were recruited from 4 separate regions including Prague, České Budějovice (in Southern
837 Bohemia), as well as Brno and Olomouc in the east of the country. Each individual residence
838 was geocoded to its administrative region. All locations outside the country of recruitment were
839 labeled as “Abroad”. A multi-membership mixed model was used to account for the full list of
840 regions in which each subject resided, as well as the proportion of life spent in that region
841 before diagnosis. As dependent variable, signatures were inverse-normal transformed. Models
842 were adjusted for sex and age of diagnosis (fixed effects). The regional effect was treated as
843 random effect.

844

845 **Data availability**

846 Whole genome sequencing data and patient metadata are deposited in the European
847 Genome-phenome Archive (EGA) associated with study EGAS00001003542. All other data is
848 provided in the accompanying Supplementary Tables.

849

850 **Code availability**

851 All algorithms used for data analysis are publicly available with repositories noted within the
852 respective method sections and in the accompanying reporting summary. Code used for
853 regression analysis and figures is available at:

854 https://gitlab.com/Mutographs/Mutographs_RCC.

855

856 **Methods references**

- 857 1. Moody, S. *et al.* Mutational signatures in esophageal squamous cell carcinoma from
858 eight countries with varying incidence. *Nat Genet* **53**, 1553–1563 (2021).
- 859 2. Whalley, J. P. *et al.* Framework for quality assessment of whole genome cancer
860 sequences. *Nat Commun* **11**, 5040 (2020).

- 861 3. Bergmann, E. A., Chen, B.-J., Arora, K., Vacic, V. & Zody, M. C. Conpair:
862 concordance and contamination estimator for matched tumor–normal pairs.
863 *Bioinformatics* **32**, 3196–3198 (2016).
- 864 4. Van Loo, P. *et al.* Allele-specific copy number analysis of tumors. *Proceedings of the*
865 *National Academy of Sciences* **107**, 16910–16915 (2010).
- 866 5. Nik-Zainal, S. *et al.* The life history of 21 breast cancers. *Cell* **149**, 994–1007 (2012).
- 867 6. Jones, D. *et al.* cgpCaVEManWrapper: Simple Execution of CaVEMan in Order to
868 Detect Somatic Single Nucleotide Variants in NGS Data. *Curr Protoc Bioinformatics*
869 **56**, (2016).
- 870 7. Raine, K. M. *et al.* cgpPindel: Identifying Somatically Acquired Insertion and Deletion
871 Events from Paired End Sequencing. *Curr Protoc Bioinformatics* **52**, (2015).
- 872 8. Kim, S. *et al.* Strelka2: fast and accurate calling of germline and somatic variants. *Nat*
873 *Methods* **15**, 591–594 (2018).
- 874 9. Bergstrom, E. N. *et al.* SigProfilerMatrixGenerator: a tool for visualizing and exploring
875 patterns of small mutational events. *BMC Genomics* **20**, 685 (2019).
- 876 10. Liu, M., Wu, Y., Jiang, N., Boot, A. & Rozen, S. G. mSigHdp: hierarchical Dirichlet
877 process mixture modeling for mutational signature discovery. *bioRxiv*
878 2022.01.31.478587 (2022) doi:10.1101/2022.01.31.478587.
- 879 11. Islam, S. M. A. *et al.* Uncovering novel mutational signatures by de novo extraction
880 with SigProfilerExtractor. *Cell genomics* **2**, None (2022).
- 881 12. Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer.
882 *Nature* **578**, 94–101 (2020).
- 883 13. Senkin, S. MSA: reproducible mutational signature attribution with confidence based
884 on simulations. *BMC Bioinformatics* **22**, 540 (2021).
- 885 14. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues.
886 *Cell* **171**, 1029-1041.e21 (2017).
- 887 15. Dentro, S. C., Wedge, D. C. & van Loo, P. Principles of Reconstructing the Subclonal
888 Architecture of Cancers. *Cold Spring Harb Perspect Med* **7**, (2017).

- 889 16. Wilcoxon, F. Individual Comparisons by Ranking Methods. *Biometrics Bulletin* **1**, 80
890 (1945).
- 891 17. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and
892 Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series*
893 *B (Methodological)* **57**, 289–300 (1995).
- 894 18. Sung, H. *et al.* Global Cancer Statistics 2020: GLOBOCAN Estimates of Incidence
895 and Mortality Worldwide for 36 Cancers in 185 Countries. *CA Cancer J Clin* **71**, 209–
896 249 (2021).
- 897 19. Dušek, L. *et al.* Epidemiology of Malignant Tumours in the Czech Republic
898 [online]. Masaryk University, Czech Republic, [2005]. <http://www.svod.cz/>. Version 7.0
899 [2007], ISSN 1802 – 8861.
- 900 20. Liu, M. *et al.* Association studies of up to 1.2 million individuals yield new insights into
901 the genetic etiology of tobacco and alcohol use. *Nat Genet* **51**, 237–244 (2019).
- 902 21. Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body
903 mass index in ~700000 individuals of European ancestry. *Hum Mol Genet* **27**, 3641–
904 3649 (2018).
- 905 22. Lagou, V. *et al.* Sex-dimorphic genetic effects and novel loci for fasting glucose and
906 insulin variability. *Nat Commun* **12**, 24 (2021).
- 907 23. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
908 *Nature* **562**, 203–209 (2018).
- 909 24. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry
910 in unrelated individuals. *Genome Res* **19**, 1655–1664 (2009).
- 911 25. Shim, H. *et al.* A Multivariate Genome-Wide Association Analysis of 10 LDL
912 Subfractions, and Their Response to Statin Treatment, in 1868 Caucasians. *PLoS*
913 *One* **10**, e0120758 (2015).
- 914 26. Choi, S. W. & O'Reilly, P. F. PRSice-2: Polygenic Risk Score software for biobank-
915 scale data. *Gigascience* **8**, (2019).

- 916 27. Lofffield, E. *et al.* Novel Biomarkers of Habitual Alcohol Intake and Associations With
917 Risk of Pancreatic and Liver Cancers and Liver Disease Mortality. *JNCI: Journal of the*
918 *National Cancer Institute* **113**, 1542–1550 (2021).
- 919 28. Shearer, J. J. *et al.* Serum Concentrations of Per- and Polyfluoroalkyl Substances and
920 Risk of Renal Cell Carcinoma. *J Natl Cancer Inst* **113**, 580–587 (2021).
- 921 29. Gao, J., Meyer, K., Borucki, K. & Ueland, P. M. Multiplex Immuno-MALDI-TOF MS for
922 Targeted Quantification of Protein Biomarkers and Their Proteoforms Related to
923 Inflammation and Renal Dysfunction. *Anal Chem* **90**, 3366–3373 (2018).

924

925 **Acknowledgements**

926 The authors would like to thank Laura O'Neill, Kirsty Roberts, Katie Smith, Maisie Farenden,
927 Siobhan Austin-Guest and the staff of DNA Pipelines at the Wellcome Sanger Institute for their
928 contribution. We are grateful for the support provided by Maja Milosevic, Christophe
929 Lallemand, Helene Renard, Aude Bardot, Andreea Spanu and Nivonirina Robinot as well as
930 IARC General Services, including the Laboratory Services and Biobank team led by Zisis
931 Kozlakidis, the Section of Support to Research overseen by Tamas Landesz and the Evidence
932 Synthesis and Classification Section led by Ian Cree, under IARC regular budget funding. The
933 authors would like to thank Gislaine Bergo, Riley Cox and Juliana Oliveira for help with
934 data/sample preparation and processing. The authors would also like to acknowledge the
935 contributions of the Leeds Biobanking and Sample Processing Lab, the Leeds Multidisciplinary
936 RTB and the Leeds NIHR BioRTB for provision of samples. The authors would like to thank
937 Peter Campbell, Inigo Martincorena, Tim Butler, Daniela Mariosa, Laura Torrens Fontanals,
938 Wellington Oliveira Dos Santos, Hana Zahed, Marc Gunter, Maggie Blanks and Mimi McCord
939 for useful discussions. The authors would also like to thank all the patients and their families
940 involved in this study.

941

942

943

944 **Funding**

945 This work was delivered as part of the Mutographs team supported by the Cancer Grand
946 Challenges partnership funded by Cancer Research UK (C98/A24032). This work was
947 supported by the Wellcome Trust grants 206194 and 220540/Z/20/A. The work was also partly
948 funded by Barretos Cancer Hospital, the Public Ministry of Labor of Campinas (Research,
949 Prevention, and Education of Occupational Cancer, 2015 to R.M.R.), and by Hospital de
950 Clínicas de Porto Alegre (180330 to P.A.-P., M.B., B.S.N.). The work was also partly supported
951 by the Practical Research Project for Innovative Cancer Control from the Japan Agency for
952 Medical Research and Development (AMED) (JP20ck0106547h0001 to T.S.), and by the
953 National Cancer Center Japan Research and Development Fund (2020-A-7 to A.F.). The work
954 was also partly funded by the 1st and 2nd Faculties of Medicine, Charles University, Prague
955 (CAGEKID to I.H.; Occupation, Environment and Kidney Cancer in Central and Eastern
956 Europe to A.H.). The work was also partly supported by the Ministry of Health of the Czech
957 Republic (MH CZ – DRO (MMCI, 00209805) to L.F. and M.N.). Measurement of PFAS
958 compounds was funded by Division of Cancer Epidemiology and Genetics of the National
959 Cancer Institute (USA). Measurement of cystatin C was funded by Cancer Research UK
960 (C18281/A29019).

961

962 **Contributions**

963 The study was conceived, designed and supervised by M.R.S., P.B. and L.B.A. Analysis of
964 data was performed by S.Senkin, S.Moody, M.D.-G., T.C., A.F.-I., J.W., S.F., M.K., R.V.,
965 A.P.L., E.N.B., A.K., B.O., S.C., E.T., J.A., K.S.-B., R.C.C.P., V.G., D.J., J.W.T. and J.M.
966 Pathology review was carried out by B.A.-A., S.F. and M.A. Sample manipulation was carried
967 out by C.L., C.C. and P.C. Patient and sample recruitment was led or facilitated by
968 S.Sangkhathat, W.A., B.S., S.J., R.S., D.M., V.Jinga, S.R., S.Milosavljevic, M.M., S.Savic,
969 J.M.S.B, M.A., L.P., P.A.-P., M.B., B.S.N., S.M.B., M.P.C., S.C.Z., R.M.R., E.F., N.S.M.,
970 R.S.F., R.B., N.V., D.Z., A.M., O.S., V.M., L.F., M.N., I.H., A.H., V.Janout, S.C. and C.L., M.P.
971 P.K.-R., S.C., M.P., P.M.U. and M.J. contributed to data generation. Patient and sample

972 recruitment for Japanese cases was led by T.S. and A.F. Scientific project management was
973 carried out by L.H., E.C., G.S., A.C.D.C., A.F.-I. and S.P. S.Moody and S.Senkin jointly
974 contributed and were responsible for overall scientific coordination. The manuscript was
975 written by S.Senkin, S.Moody, M.R.S. and P.B. with contributions from all other authors.

976

977 **Competing interests**

978 LBA is a compensated consultant and has equity interest in io9, LLC and Genome Insight. His
979 spouse is an employee of Biotheranostics, Inc. LBA is also an inventor of a US Patent
980 10,776,718 for source identification by non-negative matrix factorization. ENB and LBA declare
981 U.S. provisional applications with serial numbers: 63/289,601; 63/269,033; and 63/483,237.
982 LBA also declares U.S. provisional applications with serial numbers: 63/366,392; 63/367,846;
983 63/412,835; and 63/492,348. VM received honoraria from Ipsen, Bayer, AstraZeneca,
984 Janssen, Astellas Pharm and MSD, and provided expert testimony to BMS, Bayer, MSD and
985 Janssen. No other authors declare any competing interests.

986

987 **Disclaimer**

988 Where authors are identified as personnel of the International Agency for Research on Cancer
989 / World Health Organization, the authors alone are responsible for the views expressed in this
990 article and they do not necessarily represent the decisions, policy or views of the International
991 Agency for Research on Cancer / World Health Organization.

992

993 **Corresponding author**

994 Correspondence to Paul Brennan.

995

996 **EXTENDED DATA FIGURE AND TABLE LEGENDS**

997 **Extended Data Fig. 1: Mutation burdens in clear cell renal cell carcinomas across**
998 **countries.**

999 Mutation burdens for single base substitutions (SBS) **(a)**, doublet base substitutions (DBS)
1000 **(b)** and small insertions and deletions (ID) **(c)** show significant differences between countries
1001 using the Kruskal-Wallis (two-sided) test (n=961 biologically independent samples). Four
1002 SBS hypermutators and four ID hypermutators above mutation burden of 30000 and 3000,
1003 respectively, were removed for clarity. Box and whiskers plots are in the style of Tukey. The
1004 line within the box is plotted at the median while the upper and lower ends are indicated 25th
1005 and 75th percentiles. Whiskers show 1.5*IQR (interquartile range) and values outside it are
1006 shown as individual data points.

1007

1008 **Extended Data Fig. 2: Principal component analysis of relative mutation counts.**

1009 PCA performed on relative mutation counts of all ccRCC tumors incorporating the six
1010 mutation classes (C>A, C>G, C>T, T>A, T>C, T>G). Principal component 1 (PC1) clearly
1011 separates the cluster of mostly Romanian cases that are enriched with AA signatures, often
1012 at high mutation burdens. Principal component 3 (PC3) identifies a cluster of mostly
1013 Japanese cases, enriched with signature SBS12.

1014

1015 **Extended Data Fig. 3: Doublet-base substitution signatures operative in clear cell**
1016 **renal cell carcinomas.**

1017 **(a)** Tumour mutation burden (TMB) plot showing the frequency and mutations per Mb for
1018 each of the decomposed DBS signatures. **(b)** Average relative attribution for doublet-base
1019 substitution (DBS) signatures across countries. Signatures contributing less than 5% on
1020 average are grouped in the 'Other' category, apart from signature DBS_D. Category named
1021 '<95% confidence' accounts for the proportion of mutation burden which could not be
1022 assigned to any signature with confidence level of at least 95%. **(c)** Decomposed DBS
1023 signatures, including reference COSMIC signatures as well as *de novo* signatures not
1024 decomposed into COSMIC reference signatures.

1025

1026 **Extended Data Fig. 4: Small insertions and deletion signatures operative in clear cell**
1027 **renal cell carcinomas.**

1028 **(a)** Tumour mutation burden (TMB) plot showing the frequency and mutations per Mb for
1029 each of the decomposed ID signatures. **(b)** Average relative attribution for small insertion and
1030 deletion (ID) signatures across countries. Signatures contributing less than 5% on average
1031 are grouped in the 'Others' category, apart from signature ID_C. Category named '<95%
1032 confidence' accounts for the proportion of mutation burden which could not be assigned to
1033 any signature with confidence level of at least 95%. **(c)** Decomposed ID signatures, including
1034 reference COSMIC signatures as well as *de novo* signatures not decomposed into COSMIC
1035 reference signatures.

1036

1037 **Extended Data Fig. 5: Correlation amongst signatures SBS22a, SBS22b, DBS_D, ID_C.**

1038

1039 **Extended Data Table 1: Presence of signatures SBS22a, SBS22b, DBS_D, ID_C across**
1040 **countries.**

1041

1042 **Extended Data Fig. 6: Single base substitution signatures showing significant**
1043 **differences in attributed mutation burden between countries.**

1044 Signatures SBS40a **(a)** and SBS40b **(b)** were more prevalent in high-incidence regions of
1045 Czech Republic and Lithuania. Signatures SBS22a **(c)** and SBS22b **(d)** were enriched in
1046 Romania and Serbia. SBS1 **(e)**, SBS5 **(f)** and SBS4 **(g)** showed moderate differences across
1047 countries. Signature SBS12 **(h)** is highly prevalent in Japan. Five SBS1 hypermutators above
1048 mutation burden of 1000 were removed for clarity. Box and whiskers plots are in the style of
1049 Tukey. The line within the box is plotted at the median while the upper and lower ends are
1050 indicated 25th and 75th percentiles. Whiskers show 1.5*IQR (interquartile range) and values
1051 outside it are shown as individual data points.

1052

1053

1054 **Extended Data Fig. 7: Association of mutational signatures with incidence of renal**
1055 **cancer.**

1056 Number of mutations attributed to signatures **(a)** SBS40a, **(b)** ID5 and **(c)** ID8 against age-
1057 standardized incidence rate (ASR) of kidney cancer in each of the eleven countries
1058 represented in the cohort. The p-values shown are for ASR variable in linear regressions
1059 across all cases, adjusted for sex and age of diagnosis.

1060

1061 **Extended Data Fig. 8: Association of mutation burden with incidence of renal cancer.**

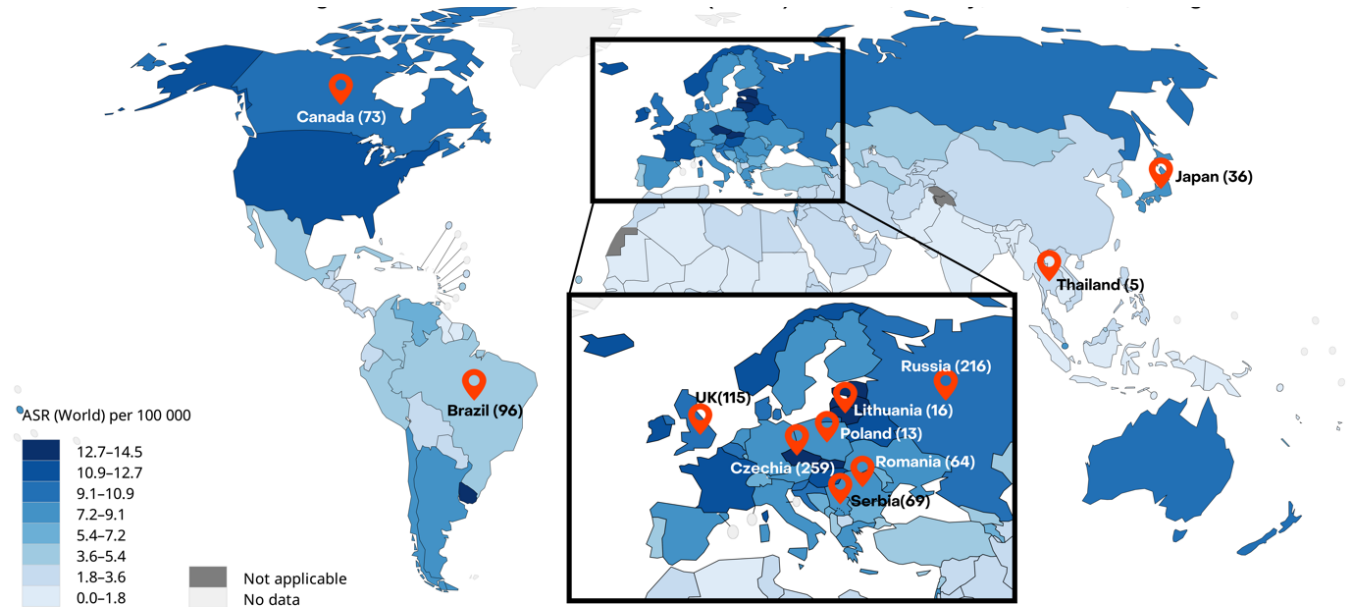
1062 Association of age-standardized rates (ASR) of kidney cancer incidence with SBS **(a)**, DBS
1063 **(b)** and ID **(c)** mutation burdens across countries. Romania (ASR \approx 7.7) and Serbia (ASR \approx 7.4)
1064 were removed due to the region-specific exposure to Aristolochic acid, with AA-related
1065 signatures accounting for a large proportion of mutation burden in these countries. The p-
1066 values shown are for ASR variable in linear regressions across all cases, adjusted for sex
1067 and age of diagnosis.

1068

1069 **Extended Data Fig. 9: Evolutionary analysis of mutational signatures in ccRCC.**

1070 Comparison of mutational signatures between clonal and subclonal mutations. Lines show
1071 the change in relative activity between the clonal mutations (main) and subclonal mutations
1072 (sub) within a sample. Blue and red lines represent an activity change of more than 6% (blue
1073 indicates higher in the clonal mutations; red indicates higher in the subclonal mutations). Bar
1074 plots show the distribution of activities in samples where the signature was present in the
1075 clonal and/or subclonal mutations; this number is represented in the title of each plot as
1076 X/223 for each signature. Black bars indicate one standard deviation away from the mean.
1077 Significance was assessed using a two-sided Wilcoxon signed-rank test, and q-values were
1078 generated using the Benjamini-Hochberg Procedure.

Fig.1



All rights reserved. The designations employed and the presentation of the material in this publication do not imply the expression of any opinion whatsoever on the part of the World Health Organization / International Agency for Research on Cancer concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. Dotted and dashed lines on maps represent approximate borderlines for which there may not yet be full agreement.

Data source: GLOBOCAN 2020
Map production: IARC
(<http://gco.iarc.fr/today>)
World Health Organization



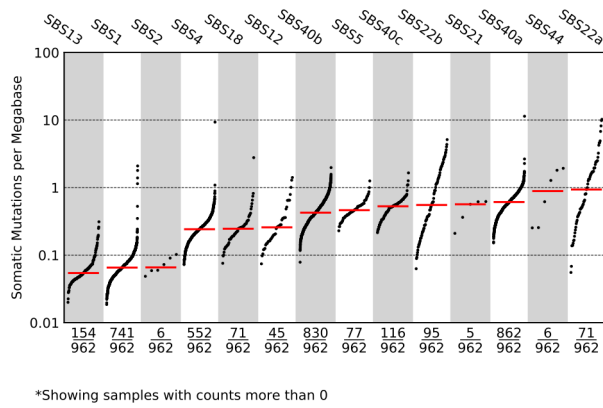
© International Agency for Research on Cancer 2020
All rights reserved

Table 1

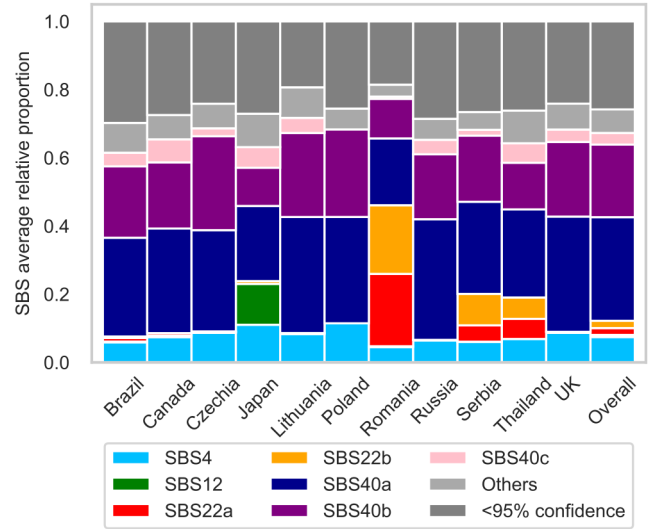
Country (ASR/100,000)	Brazil (4.5)	Canada (10.4)	Czechia (14.4)	Japan (7.6)	Lithuania (14.5)	Poland (8.1)	Romania (7.7)	Russia (10.3)	Serbia (7.4)	Thailand (1.8)	UK (10.3)	Total (4.6)	
Number of cases	96	73	259	36	16	13	64	216	69	5	115	962	
Sex	Female	44	22	93	8	9	5	25	98	30	4	42	380
	Male	52	51	166	28	7	8	39	118	39	1	73	582
Age at diagnosis (years)	0-45	15	6	27	3	1	2	6	43	16	0	6	125
	45-55	20	17	51	5	0	6	10	44	11	0	22	186
	55-65	30	17	77	8	9	1	20	91	27	2	41	323
	65-75	24	27	72	13	4	4	20	32	9	2	31	238
	75+	7	6	32	7	2	0	8	6	6	1	15	90
Year of recruitment	1999-2005			93			13	14	18			138	
	2005-2010			111				19	70	1		232	
	2010-2015		9	55	28			31	116	68		348	
	2015-2020	96	64		8	16			12		5	244	
Stage	I	28	3	123	24	6	0	33	94	32		53	396
	II	2	0	42	1	0	6	12	24	4		8	99
	III	16	23	46	6	5	5	18	65	26		38	248
	IV	7	10	38	5	2	2	1	33	7		16	121
	Missing	43	37	10		3					5		98
Body mass index	<20	3	2	5	2	0	2	2	9	8	0	6	39
	20-25	21	10	100	25	2	3	17	84	28	3	23	316
	25-30	35	24	85	7	6	6	30	40	20	1	45	299
	>30	37	37	69	2	8	2	14	83	13	1	41	307
	Missing							1					1
Hypertension	No	45	28	129	16	5	9	39	125	28	2	58	484
	Yes	51	44	130	20	10	4	24	91	41	3	56	474
	Missing		1			1		1				1	4
Diabetes	No	76	55	130	29	9		45	186	61	3	95	689
	Yes	20	16	36	7	7		4	12	8	2	20	132
	Missing		2	93			13	15	18				141
Family history of ccRCC	No	90	42	165	35	16		54	192	67	5	102	726
	Yes	5	4	22	1	0		1	6	2	0	3	43
	Missing	1	27	72			13	9	18			10	193
Tobacco status	Current smoker	23	21	66	9	4	6	11	52	18	1	28	239
	Ex-smoker	21	30	62	15	3	3	15	27	15	0	44	235
	Never	52	22	131	11	9	4	37	137	36	4	43	486
	Missing				1			1					2
PFOA (ng/mL)	Mean (st. dev.)	0.7 (0.5)	1.6 (1.1)	3.4 (2.1)		1.3 (0.6)	5.4 (4.1)	1.3 (0.9)	1.5 (1.4)	1.3 (0.6)	2.2 (2.2)	3.3 (1.7)	2.2 (1.9)

Fig. 2

a



b



c



Fig. 3

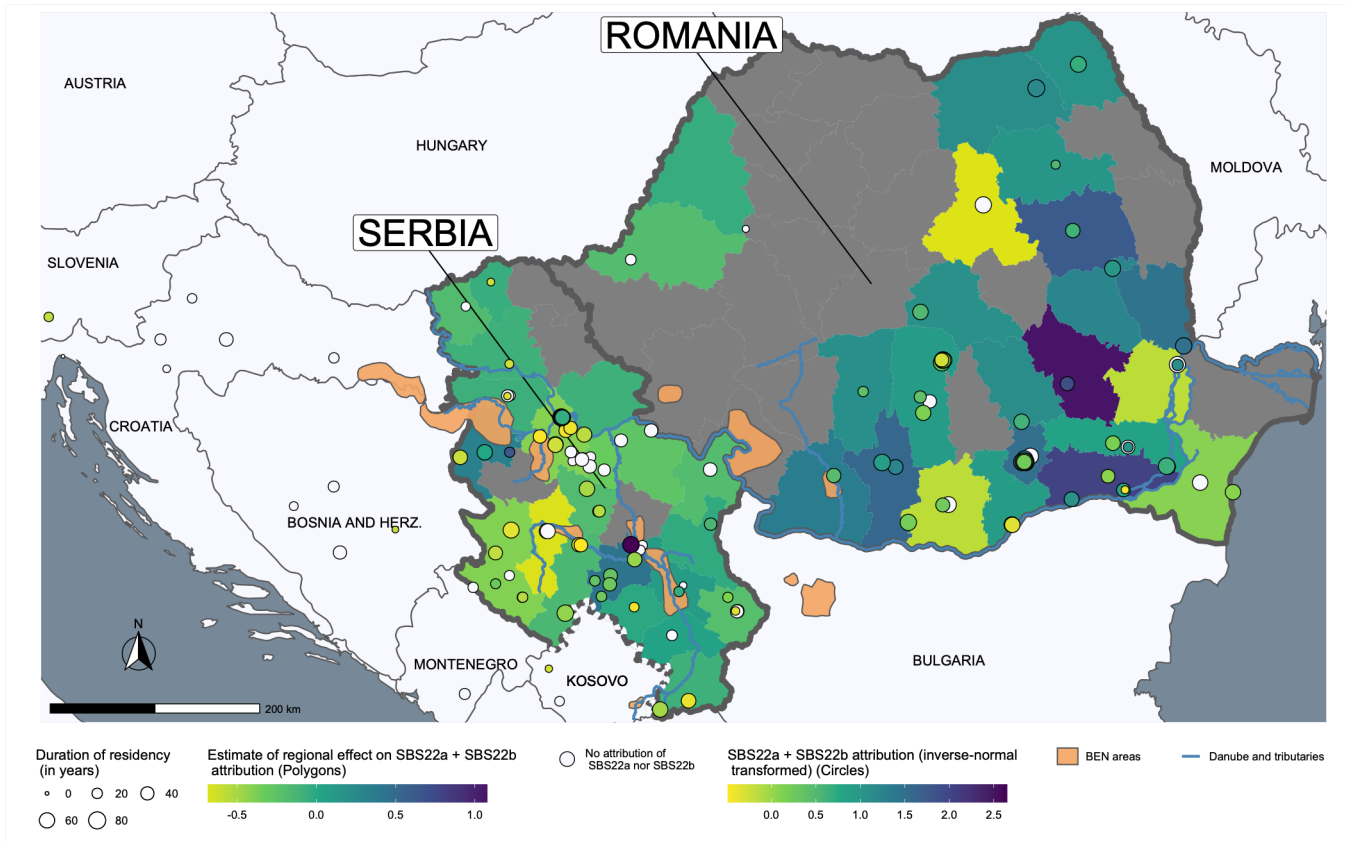
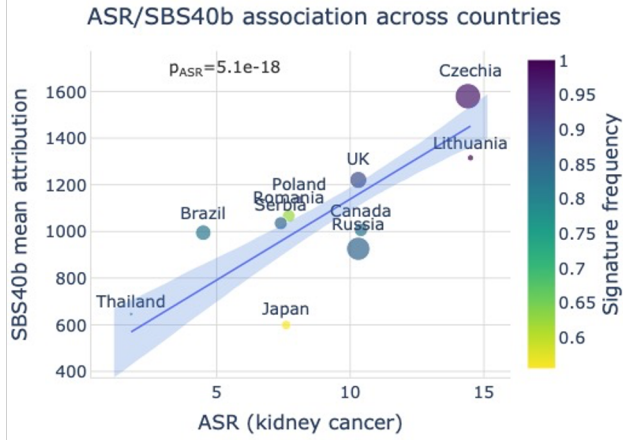
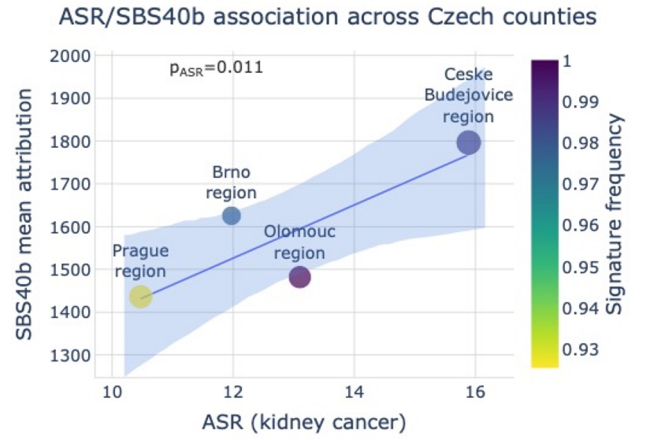


Fig. 4

a



b



c

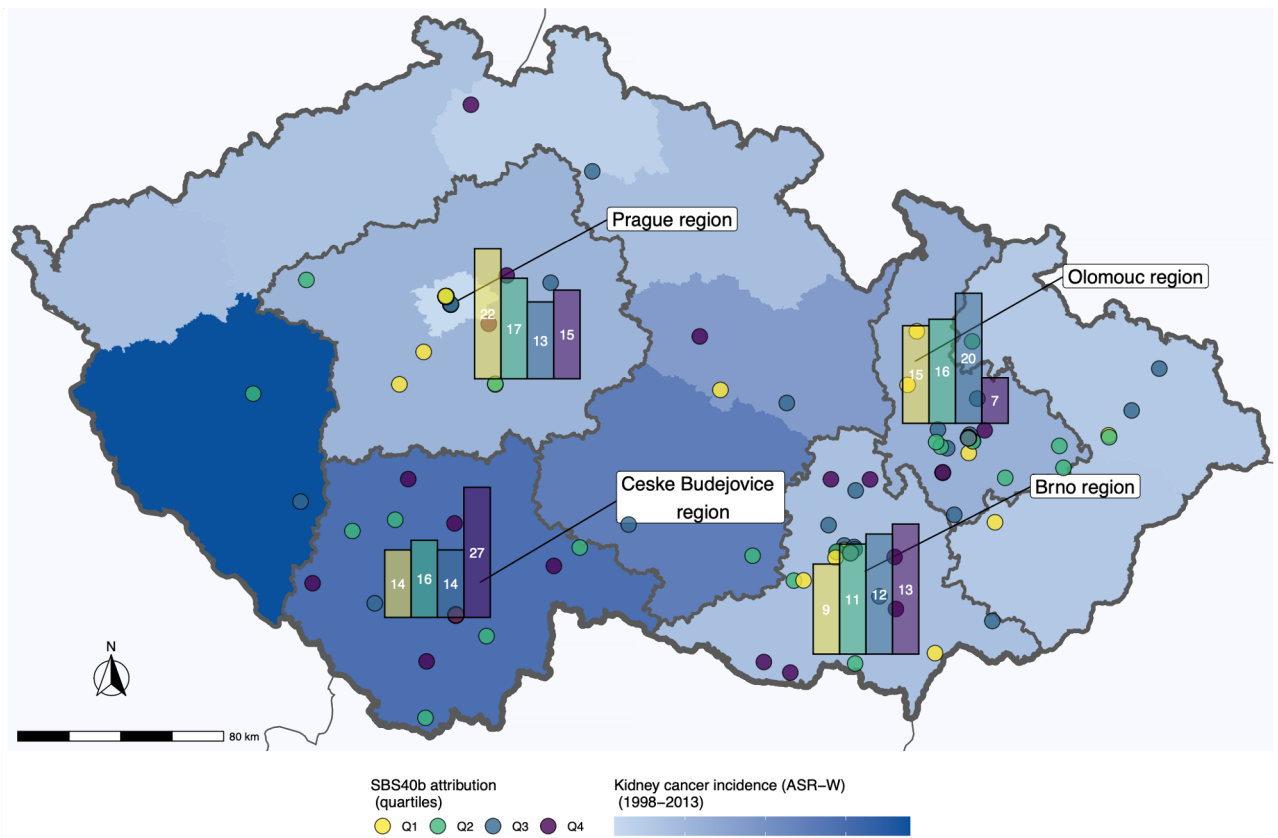
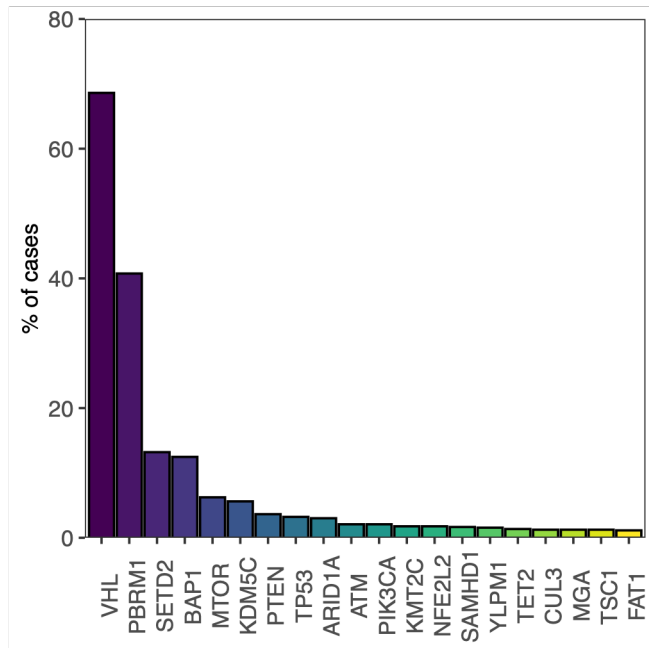
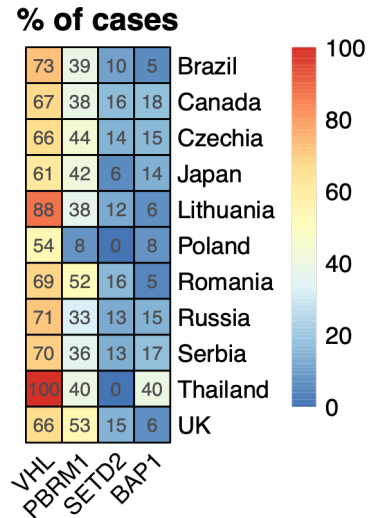


Fig. 5

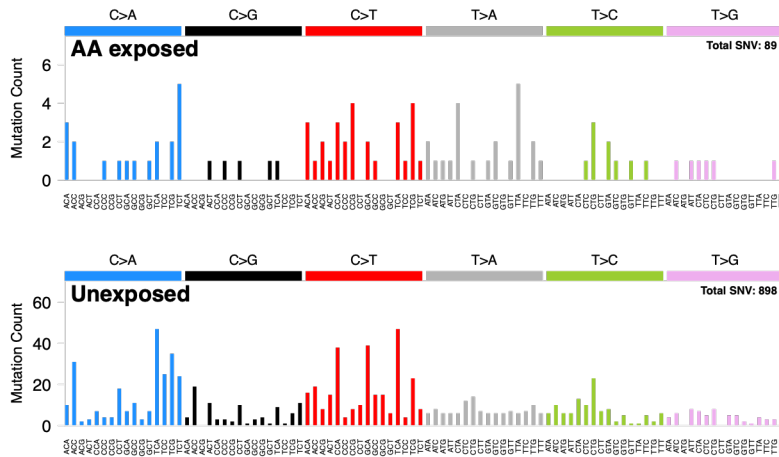
a



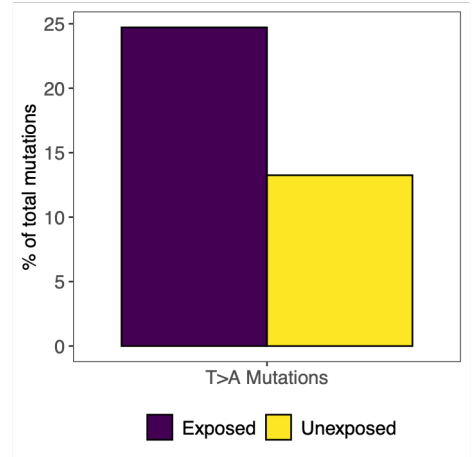
b



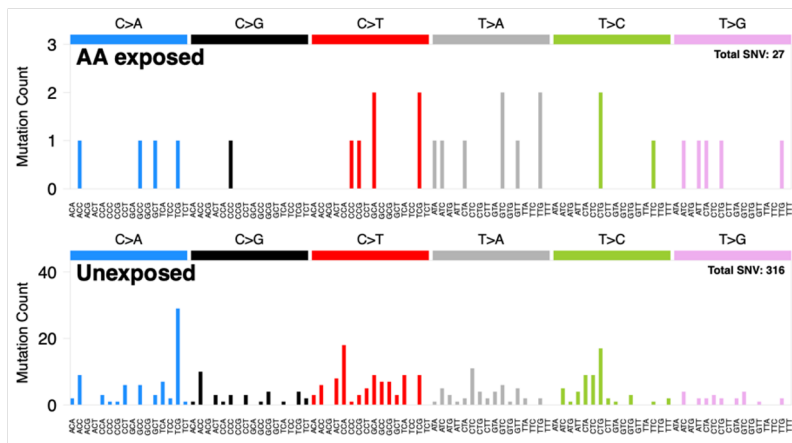
c



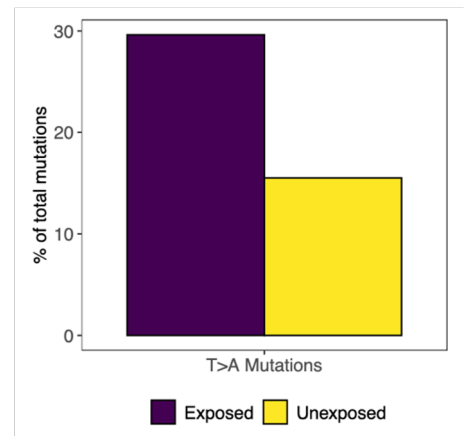
d



e

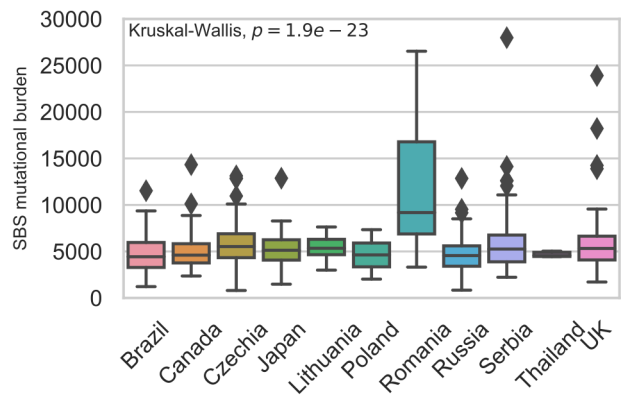


f

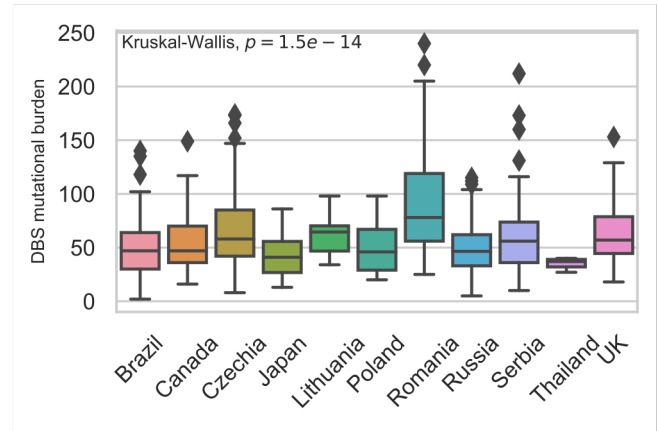


Extended Data Fig. 1.

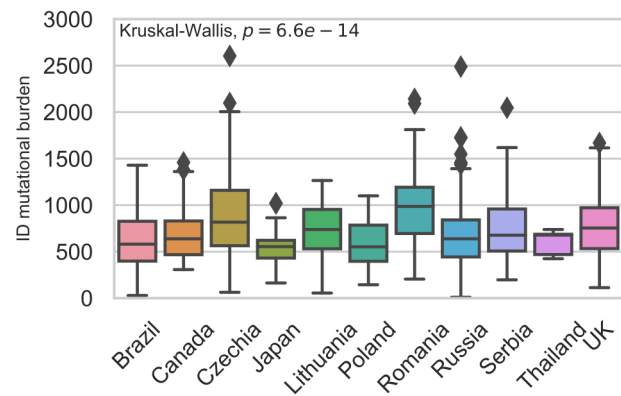
a



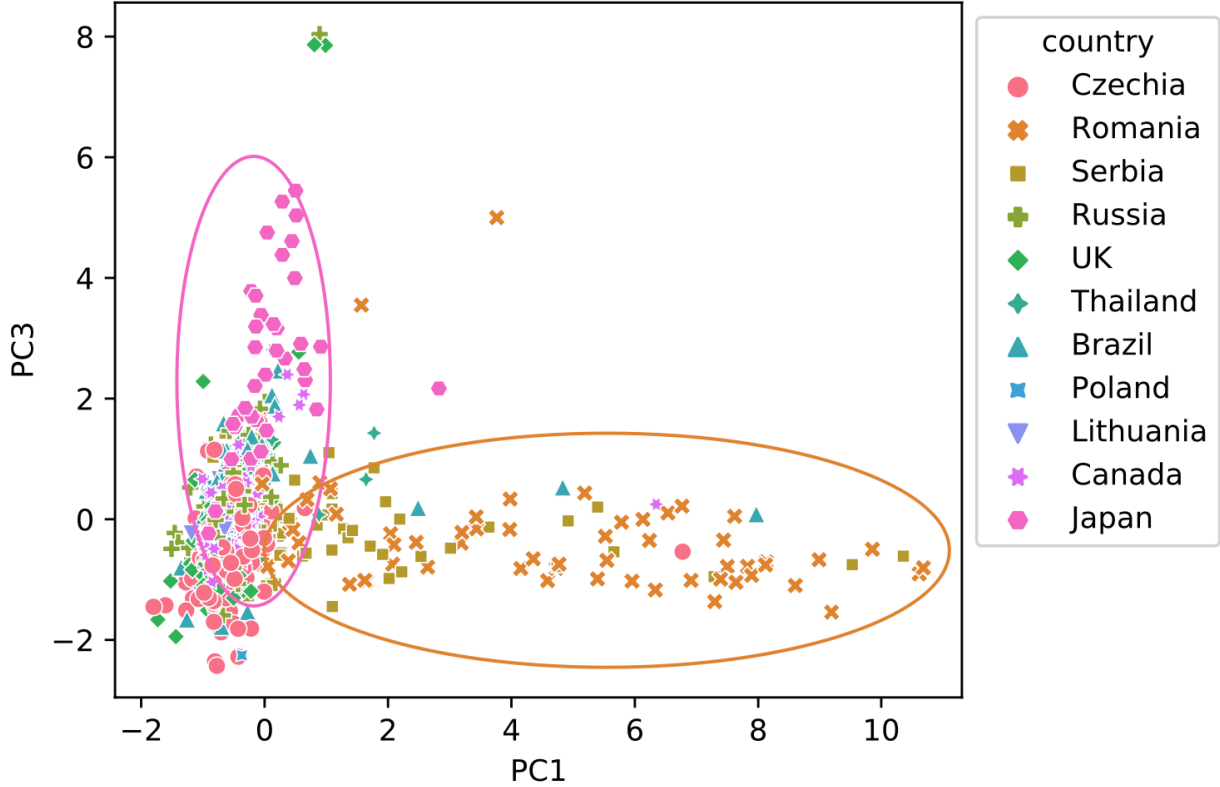
b



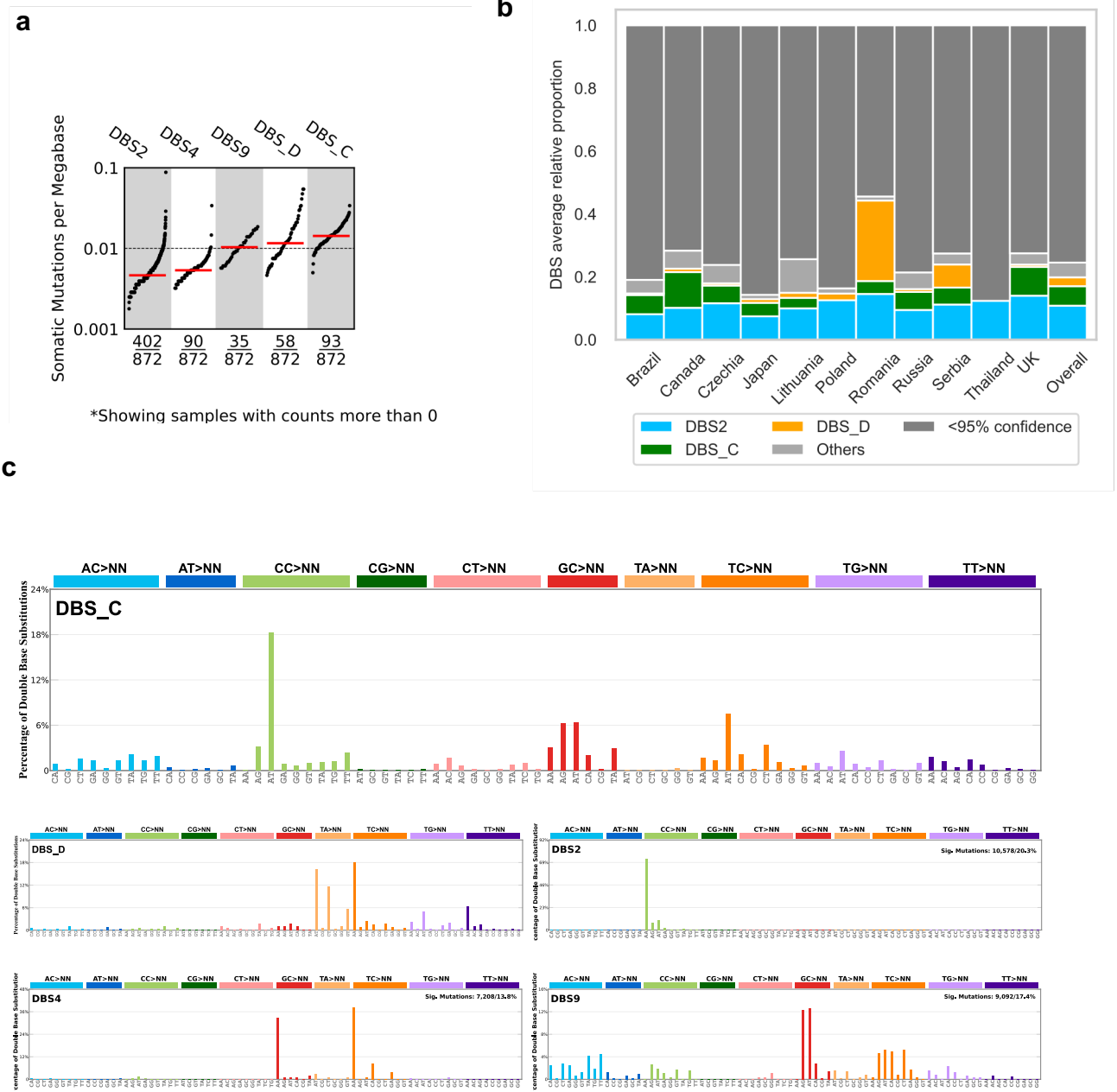
c



Extended Data Fig. 2.

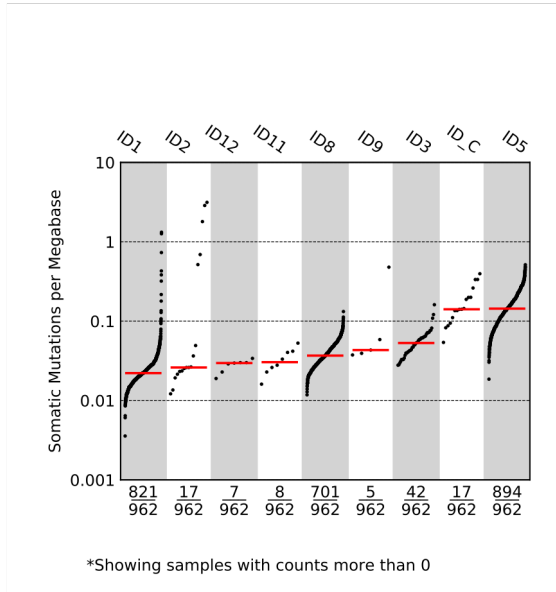


Extended Data Fig. 3.

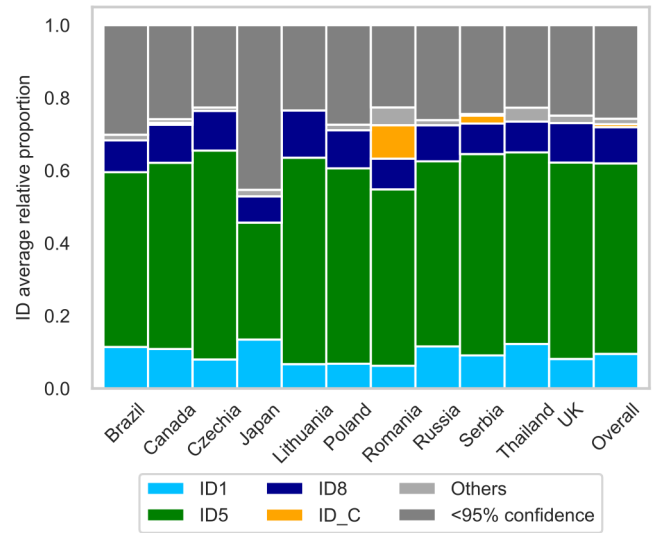


Extended Data Fig. 4.

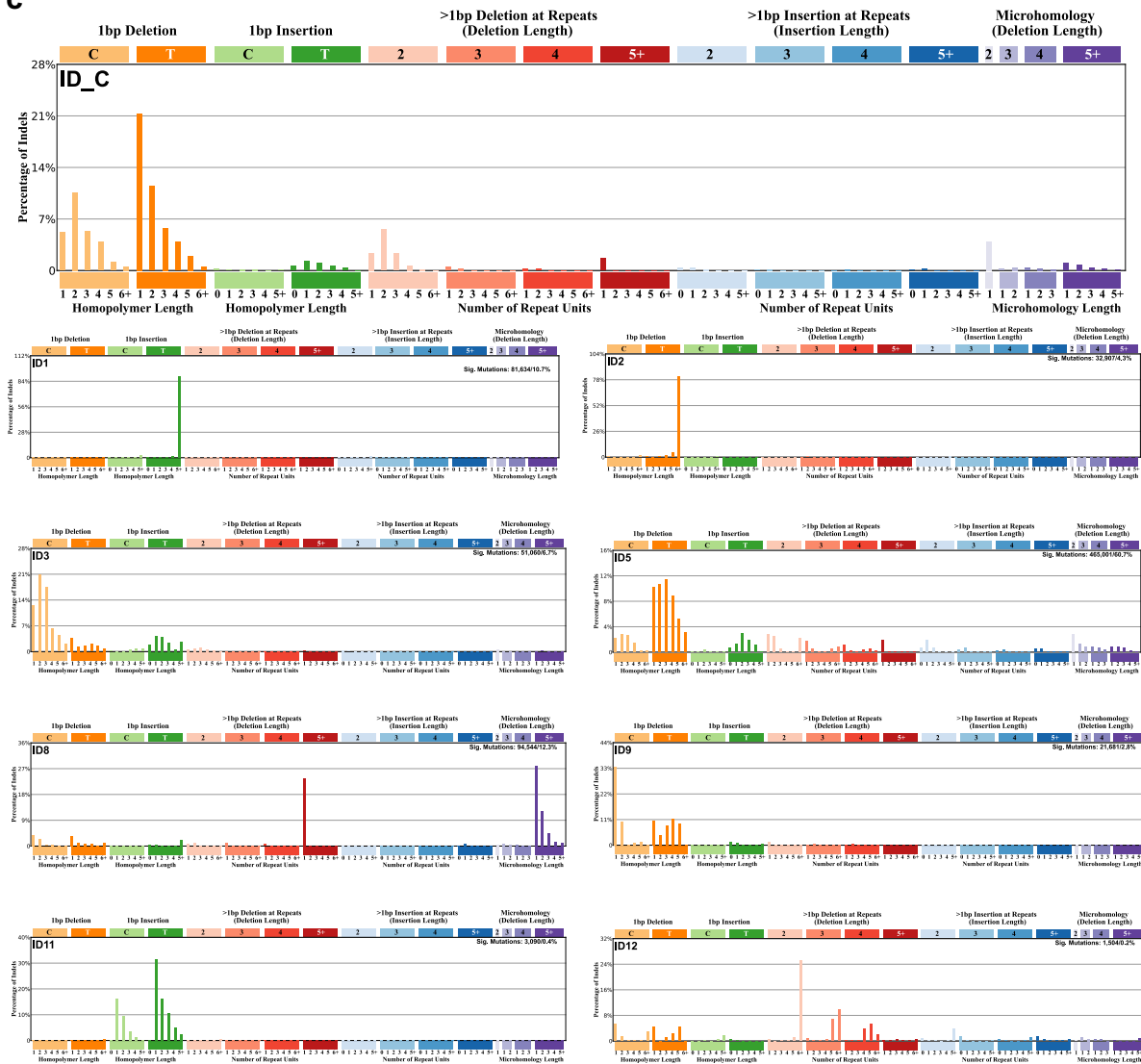
a



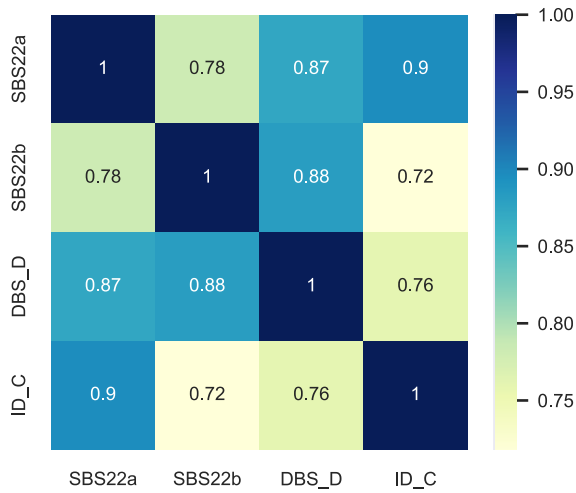
b



c



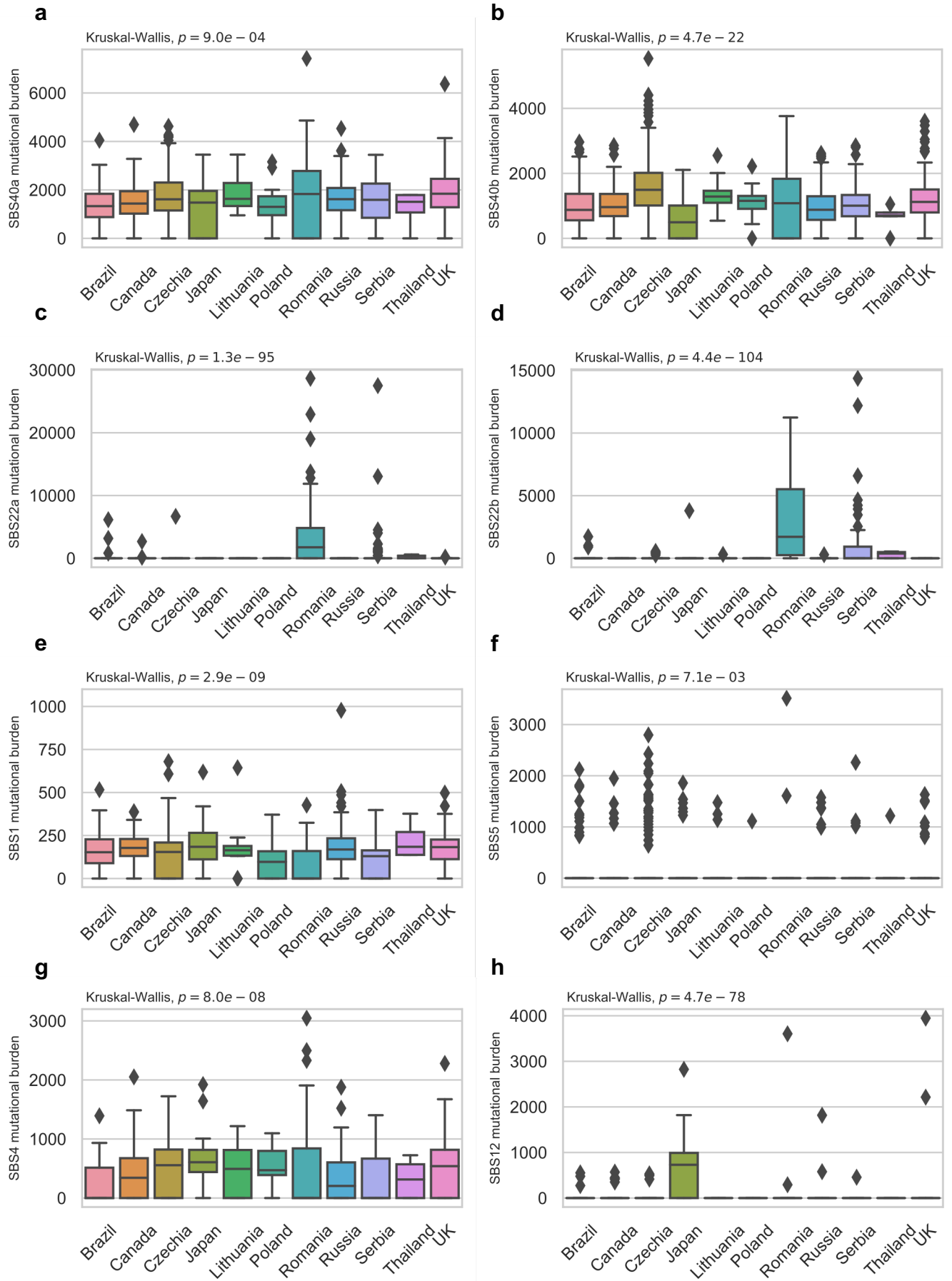
Extended Data Fig. 5.



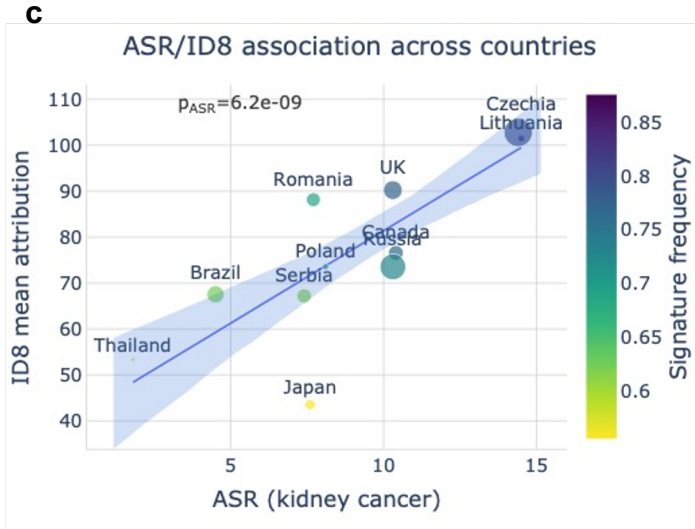
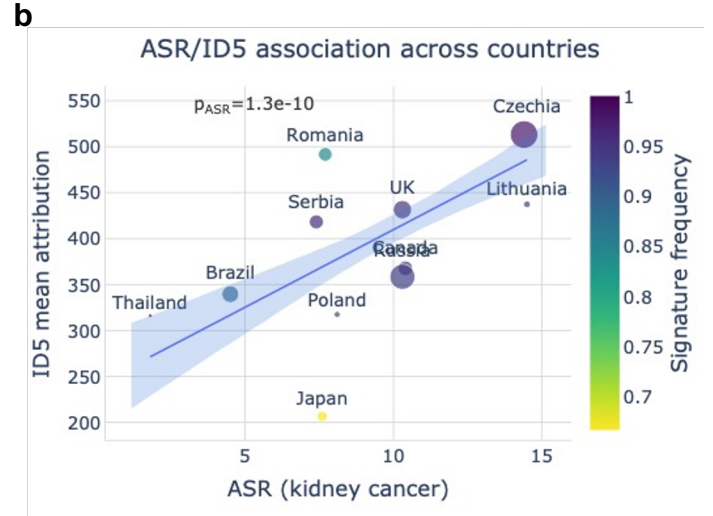
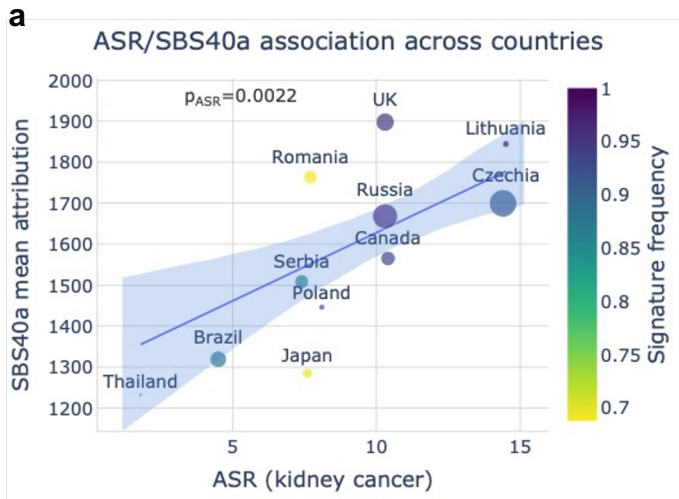
Extended Data Table 1.

Country	N cases	SBS22a (%)	SBS22b (%)	DBS_D (%)	ID_C (%)	SBS22a or SBS22b (%)	Any (%)
Romania	64	45 (70.3)	48 (75.0)	42 (65.6)	13 (20.3)	53 (82.8)	54 (84.4)
Serbia	69	16 (23.2)	33 (47.8)	11 (15.9)	3 (4.3)	35 (50.7)	36 (52.2)
Thailand	5	3 (60.0)	3 (60.0)	0 (0.0)	0 (0.0)	4 (80.0)	4 (80.0)
Brazil	96	3 (3.1)	3 (3.1)	1 (1.0)	0 (0.0)	3 (3.1)	3 (3.1)
Canada	73	2 (2.7)	0 (0.0)	2 (2.7)	1 (1.4)	2 (2.7)	3 (4.1)
Czechia	259	1 (0.4)	5 (1.9)	32 (12.4)	0 (0.0)	6 (2.3)	37 (14.3)
UK	115	1 (0.9)	0 (0.0)	31 (27.0)	0 (0.0)	1 (0.9)	31 (27.0)
Russia	216	0 (0.0)	1 (0.5)	26 (12.0)	0 (0.0)	1 (0.5)	27 (12.5)
Poland	13	0 (0.0)	0 (0.0)	1 (7.7)	0 (0.0)	0 (0.0)	1 (7.7)
Lithuania	16	0 (0.0)	1 (6.2)	1 (6.2)	0 (0.0)	1 (6.2)	2 (12.5)
Japan	36	0 (0.0)	1 (2.8)	1 (2.8)	0 (0.0)	1 (2.8)	1 (2.8)

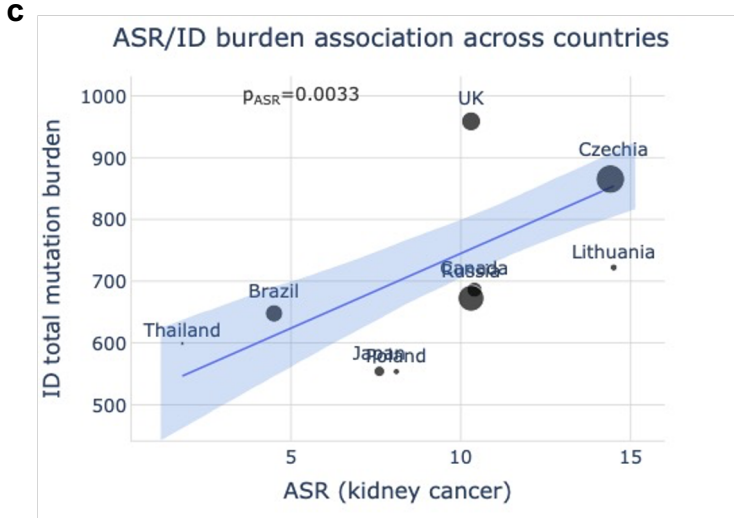
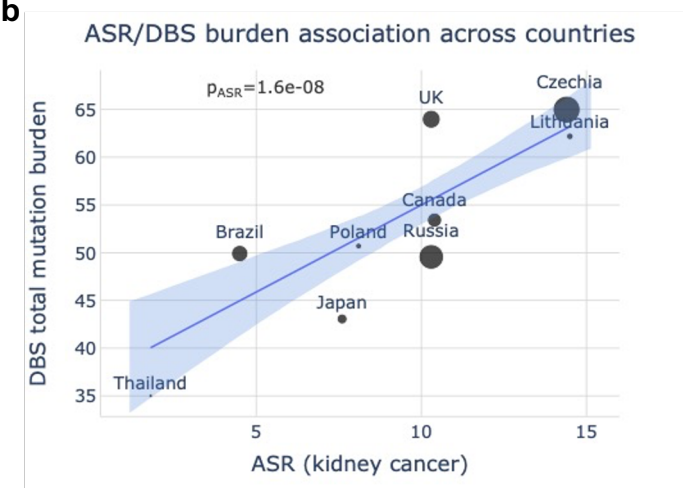
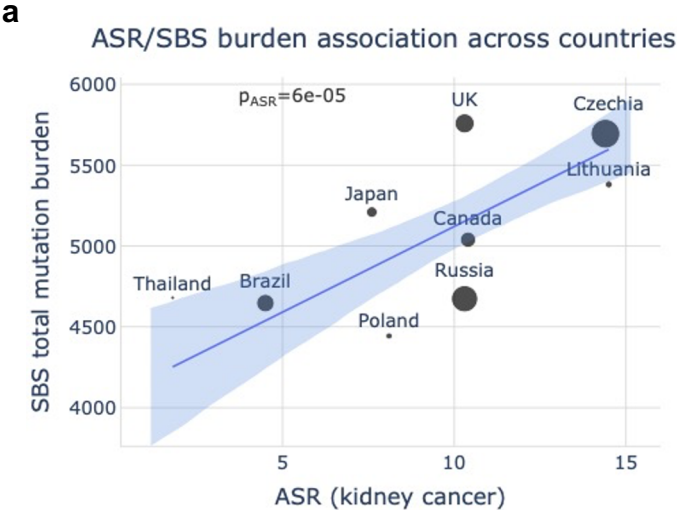
Extended Data Fig. 6.



Extended Data Fig. 7.



Extended Data Fig. 8.



Extended Data Fig. 9.

