

Open access • Posted Content • DOI:10.1101/2020.12.13.422570

Uncovering novel mutational signatures by de novo extraction with SigProfilerExtractor — Source link 🖸

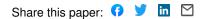
S M Ashiqul Islam, Yang Wu, Marcos Díaz-Gay, Erik N. Bergstrom ...+34 more authors

Institutions: University of California, San Diego, National University of Singapore, Nvidia, Wellcome Trust Sanger Institute ...+7 more institutions

Published on: 13 Dec 2020 - bioRxiv (Cold Spring Harbor Laboratory)

Related papers:

- The Repertoire of Mutational Signatures in Human Cancer
- · Pan-cancer analysis of whole genomes
- · Signatures of mutational processes in human cancer
- SigProfilerMatrixGenerator: a tool for visualizing and exploring patterns of small mutational events.
- · Clock-like mutational processes in human somatic cells



Uncovering novel mutational signatures by *de novo* extraction with 1

SigProfilerExtractor 2

- S M Ashiqul Islam^{1,2,3}, Yang Wu⁴, Marcos Díaz-Gay^{1,2,3}, Erik N Bergstrom^{1,2,3}, Yudou He^{1,2,3}, 3
- Mark Barnes^{1,2,3}, Mike Vella⁵, Jingwei Wang⁶, Jon W Teague⁶, Peter Clapham⁶, Sarah Moody⁶, 4
- Sergey Senkin⁷, Yun Rose Li⁸, Laura Riva⁶, Tongwu Zhang⁹, Andreas J Gruber^{10,11}, Raviteja 5
- Vangara¹², Christopher D Steele¹³, Burçak Otlu^{1,2,3}, Azhar Khandekar^{1,2,3}, Ammal Abbasi^{1,2,3}, 6
- 7 Laura Humphreys⁶, Natalia Syulyukina², Samuel W Brady¹⁴, Boian S Alexandrov¹², Nischalan
- 8 Pillay^{13,15}, Jinghui Zhang¹⁴, David J Adams⁶, Iñigo Martincorena⁶, David C Wedge^{10,11}, Maria
- 9 Teresa Landi⁹, Paul Brennan⁷, Michael R Stratton⁶, Steven G Rozen⁴, and Ludmil B
- 10 Alexandrov^{1,2,3*}
- 11

12

13 Affiliations

- 14 ¹Department of Cellular and Molecular Medicine, UC San Diego, La Jolla, CA, 92093, USA
- 15 ²Department of Bioengineering, UC San Diego, La Jolla, CA, 92093, USA
- 16 ³Moores Cancer Center, UC San Diego, La Jolla, CA, 92037, USA
- 17 ⁴Centre for Computational Biology and Programme in Cancer & Stem Cell Biology, Duke NUS
- 18 Medical School, 169857, Singapore
- 19 ⁵NVIDIA Corporation, 2788 San Tomas Expressway, Santa Clara, CA, 95051, USA
- 20 ⁶Cancer, Ageing and Somatic Mutation, Wellcome Sanger Institute, Wellcome Genome Campus,
- 21 Cambridge, CB10 1SA, UK
- 22 ⁷Genetic Epidemiology Group, International Agency for Research on Cancer, 69372 Lyon CEDEX
- 23 08. France
- 24 ⁸Helen Diller Family Comprehensive Cancer Center, San Francisco, CA, 94158, USA
- 25 ⁹Division of Cancer Epidemiology and Genetics, National Cancer Institute, Bethesda, MD, 20892, 26 USA
- 27 ¹⁰Big Data Institute, Nuffield Department of Medicine, University of Oxford, OX5 7LF, 28 UK
- 29
- ¹¹Manchester Cancer Research Centre, The University of Manchester, Manchester, M20 4GJ, UK 30 ¹²Theoretical Division, Los Alamos National Laboratory, Los Alamos, NM, 87545, USA
- 31 ¹³Research Department of Pathology, Cancer Institute, University College London, London, WC1E 32 6BT. UK
- 33 ¹⁴Department of Computational Biology, St. Jude Children's Research Hospital, Memphis, 38105,
- 34 Tennessee, USA
- 35 ¹⁵Department of Cellular and Molecular Pathology, Royal National Orthopaedic Hospital NHS
- 36 Trust, Stanmore, Middlesex, HA7 4LP, UK
- 37
- 38 *Correspondence should be addressed to L2alexandrov@health.ucsd.edu.
- 39

40 ABSTRACT

41 Mutational signature analysis is commonly performed in genomic studies surveying cancer and 42 normal somatic tissues. Here we present SigProfilerExtractor, an automated tool for accurate de 43 novo extraction of mutational signatures for all types of somatic mutations. Benchmarking with a 44 total of 33 distinct scenarios encompassing 1,900 simulated signatures operative in more than 45 60,000 unique synthetic genomes demonstrates that SigProfilerExtractor outperforms thirteen 46 other tools across all datasets with and without noise. For simulations with 5% noise, reflecting 47 high-quality genomic datasets, SigProfilerExtractor outperforms other approaches by elucidating 48 between 20% and 50% more true positive signatures while yielding more than 5-fold less false 49 positive signatures. Applying SigProfilerExtractor to 2,778 whole-genome sequenced cancers 50 reveals three previously missed mutational signatures. Two of the signatures are confirmed in 51 independent cohorts with one of these signatures associating with tobacco smoking. In summary, 52 this report provides a reference tool for analysis of mutational signatures, a comprehensive 53 benchmarking of bioinformatics tools for extracting mutational signatures, and several novel 54 mutational signatures including a signature putatively attributed to direct tobacco smoking 55 mutagenesis in bladder cancer and in normal bladder epithelium.

57 INTRODUCTION

75

De novo extraction of mutational signatures¹ is an unsupervised machine learning approach 58 59 where a matrix, M, which corresponds to the somatic mutations in a set of cancer samples under 60 a mutational classification², is approximated by the product of two low-rank matrices, S and A. 61 The matrix **S** reflects the set of mutational signatures while the matrix **A** encompasses the 62 activities of the signatures; an activity corresponds to the number of mutations contributed by a 63 signature in a cancer sample. Algorithmically, *de novo* extraction of mutational signatures has relied on nonnegative matrix factorization (NMF)³ or on approaches mathematically analogous 64 to NMF⁴⁻⁶. The main advantage of NMF over other factorization approaches is its ability to yield 65 66 nonnegative factors that are part of the original data, thus, allowing interpretation of the 67 identified nonnegative factors³. Biologically, mutational signatures extracted from cancer 68 genomes have been attributed to exposures to environmental carcinogens, failure of DNA repair 69 pathways, infidelity/deficiency of replicating polymerases, iatrogenic events, and others⁷⁻¹⁴. 70 71 Since we introduced the mathematical concept of mutational signatures¹, a number of 72 computational frameworks have been developed for performing *de novo* extraction of mutational 73 signatures (**Table 1**)¹⁴⁻²⁷. Notably, the majority of existing *de novo* extraction tools (*i*) 74 predominately support the simplest mutational classification, viz., SBS-96 which encompasses

76 selection for the number of mutational signatures; *(iii)* do not identify a robust solution leading to

single base substitutions with their immediate 5' and 3' sequence context²; (ii) lack automatic

- different results following re-analysis of the same dataset; *(iv)* require pre-selection of a large
- number of priors and/or hyperparameters; (v) do not decompose *de novo* signatures to the set of

79	reference COSMIC signatures ¹⁴ . Importantly, there has been no extensive benchmark of the
80	existing tools for <i>de novo</i> extraction leading to uncertainty regarding their performance.
81	

82 To address these limitations, here, we present SigProfilerExtractor - a reference tool for *de novo* 83 extraction of mutational signatures. SigProfilerExtractor allows analysis of all types of 84 mutational classifications, performs automatic selection of the number of signatures, yields 85 robust solutions, requires only minimum setup, and decomposes *de novo* extracted signatures to 86 known COSMIC signatures. A comprehensive benchmark including 3,448 unique matrix 87 decompositions with SigProfilerExtractor and thirteen other tools across a total of 33 distinct 88 scenarios reveals that SigProfilerExtractor is robust to noise and that it outperforms all other 89 computational tools for de novo extraction of mutational signatures (Supplementary Tables 1-90 3). Applying SigProfilerExtractor to the recently published set of 2,778 whole-genome sequenced cancers from the Pan-Cancer Analysis of Whole Genomes (PCAWG) project²⁸ 91 92 elucidates three novel signatures that were not found in the original PCAWG analysis of 93 mutational signatures¹⁴. Two of the signatures are confirmed in independent cohorts and a 94 putative etiology of tobacco-associated mutagenesis is attributed to one of these signatures. 95

96 **RESULTS**

97 Overview of SigProfilerExtractor and its implementation

- 98 SigProfilerExtractor is implemented as a Python package, with an R wrapper, allowing users to
- 99 run it in both Python and R environments:
- 100 <u>https://github.com/AlexandrovLab/SigProfilerExtractor</u>. The tool is also extensively
- 101 documented including a detailed Wiki page: <u>https://osf.io/t6j7u/wiki/home/</u>. By default, the tool
- 102 requires only a single parameter the input dataset containing the mutational catalogues of
- 103 interest. SigProfilerExtractor supports most used formats outputted by variant calling algorithms
- 104 (e.g., VCF, MAF, etc.), which are internally converted to a matrix, *M*, by
- 105 SigProfilerMatrixGenerator². SigProfilerExtractor can also be applied to a text file containing a

106 matrix, *M*, thus supporting nonnegative matrix factorization for any custom matrix dataset. By

- 107 default, the tool decomposes the matrix *M* searching for an optimal solution for the number of
- 108 operative signatures, *k*, between 1 and 40 mutational signatures (Figure 1*a*). For each
- 109 decomposition, SigProfilerExtractor performs 500 independent factorizations and, for each
- 110 repetition, the matrix *M* is first Poisson resampled and normalized and, subsequently, factorized
- 111 with the multiplicative update NMF algorithm³ by minimizing an objective function based on the
- 112 Kullback–Leibler divergence measure²⁹ (Figure 1*b*). Custom partition clustering, that utilizes the
- 113 Hungarian algorithm³⁰ for comparing different repetitions, is applied to the 500 factorizations to

114 identify stable solutions³¹ (Figure 1*b*). Specifically, SigProfilerExtractor selects the centroids of

- 115 stable clusters as optimal solutions, thus, making these solutions resistant to fluctuations in the
- 116 input data and to the lack of uniqueness of NMF due to the potential existence of multiple
- 117 convergent stationary points in the solution³². Lastly, when applicable, the optimal set of *de novo*
- signatures are matched to the set of reference COSMIC mutational signatures (Figure 1c) with

any *de novo* signature reported as novel when it cannot be decomposed by a combination ofknown COSMIC signatures.

121

122 Framework for benchmarking tools for de novo extraction of mutational signatures

123 To allow comprehensive benchmarking of tools for *de novo* extraction of mutational signatures, 124 more than 60,000 unique synthetic cancer genomes were generated with known ground-truth 125 mutational signatures (Supplementary Note 1). These synthetic data included 32 distinct 126 noiseless scenarios and one scenario with five different levels of noise. Each scenario contained 127 between 3 and 40 known signatures operative in 200 to 3,000 simulated cancer genomes 128 (Supplementary Tables 1–3). Some of the scenarios were generated up to 20 times to account 129 for variability in the simulated data. Most noiseless scenarios (20/32) were based on SBS-96 130 mutational classification; 12 scenarios based on extended mutational classifications, *i.e.*, matrices 131 with more than 96 mutational channels, were also included (Supplementary Table 3). To avoid 132 bias in evaluating each tool's performance, three sets of SBS-96 mutational signatures were used 133 for generating the synthetic data: (i) COSMICv3 reference signatures¹⁴; (ii) SA signatures previously extracted by SignatureAnalyzer¹⁴; and *(iii)* randomly generated signatures. For 134 135 presentation simplicity, scenarios were labeled based on their complexity as easy, medium, or 136 hard. Easy scenarios were generated using ≤ 5 signatures and provide a good indication of each 137 tool's performance on approximately 7.4% of human cancer types (e.g., pediatric brain tumors). 138 Medium scenarios contained 11 to 21 signatures and biologically reflect 15.9% of cancer types 139 (e.g., cervical cancer). Hard scenarios have more than 25 signatures and reflect 59.5% of human 140 cancer types (e.g., breast, lung, liver, etc.) as well as pan-cancer datasets. In addition to the 32

- noiseless scenarios, one SBS-96 scenario with five different levels of noise, ranging between 0%
 and 10%, was included in the benchmark (Supplementary Note 1).
- 143
- 144 To compare the performance between different tools for *de novo* extraction of mutational
- signatures, we developed a standard set of evaluation metrics (Supplementary Figure 1).
- 146 Specifically, each *de novo* extracted signature is classified as either a *true positive* (TP), *false*
- 147 positive (FP), or false negative (FN) signature. An extracted signature is considered TP if it
- 148 matches one of the ground-truth signatures above a cosine similarity threshold of 0.90. In

149 contrast, a signature is classified as FP when it has a maximum cosine similarity below 0.90 with

all ground-truth signatures. Lastly, FN signatures are ground-truth signatures that were not

151 detected in the data. These standard metrics allow calculating each tool's precision, sensitivity,

152 and F₁ score. Precision is defined as $\frac{TP}{TP+FP}$, sensitivity as $\frac{TP}{TP+FN}$, and F₁ score is the harmonic

153 mean of the precision and sensitivity: $2 * \frac{Precision * Sensitivity}{Precision + Sensitivity}$

154

155 Benchmarking SigProfilerExtractor and thirteen other tools using SBS-96 noiseless data 156 SigProfilerExtractor and thirteen other tools (**Table 1**) were first applied to all noiseless 157 scenarios based on the SBS-96 mutational classification. The thirteen tools include 158 SignatureAnalyzer and SigProfiler PCAWG, a legacy MATLAB/Python version of 159 SigProfilerExtractor, which were jointly used in the PCAWG analysis of mutational signatures 160 and the derivation of the COSMICv3 set of reference mutational signatures¹⁴. Except for 161 MutSignatures which can only decompose a matrix for a fix number of signatures, all other tools 162 were applied to each scenario by using their suggested methods for selecting the number of 163 operative signatures. Apart from SignatureAnalyzer which lacks this capability, all other tools

164 were forced to extract the known number of ground-truth signatures. Results from the suggested 165 approach reflect the expected outcome from running a tool on an unknown dataset, while results 166 from the forced approach allow understanding limitations in each tool's implementation. Our 167 evaluation reveals that most tools can successfully extract mutational signatures from easy 168 scenarios with the majority of F_1 scores between 0.90 and 1.00 (Figure 2*a*). This is perhaps 169 unsurprising as many of these tools used synthetic data with ≤ 5 signatures to evaluate their 170 performance in their respective original publications^{15-17,19-26}. In contrast, medium scenarios have 171 proven to be a challenge for most tools with only SigProfilerExtractor, SigProfiler PCAWG, and 172 SignatureAnalyzer exhibiting F_1 scores above 0.90. All tools had worst performance for the hard 173 set of scenarios with F₁ scores below 0.80; only SigProfilerExtractor had an F₁ score of almost 174 0.90 (Figure 2*a*).

175

176 To evaluate whether the type of ground-truth signatures affects the *de novo* extraction, we 177 compared the ratio of F_1 scores (rF_1) from scenarios generated using COSMIC, SA, or random 178 signatures (Figure 2b). Most tools had similar performance ($rF_1 \approx 1$) between COSMIC and 179 random signatures and worst performance with SA signatures ($rF_1 < 1$). SomaticSignatures was an 180 exception as it performed well on random signatures but had similarly suboptimal performance 181 on COSMIC and SA signatures. SigProfilerExtractor outperformed all other tools regardless of 182 whether the synthetic data were generated using COSMIC, SA, or random signatures 183 (Supplementary Table 1).

184

185 To examine the performance of *de novo* extraction between the suggested and forced selection of 186 the total number of signatures, we evaluated rF_1 across all medium and hard scenarios (**Figure**

187 2c). SigProfilerExtractor exhibited almost identical F₁ scores between the suggested and forced 188 selection indicating a good performance of the automatic selection algorithm. Most other tools 189 had similar F_1 scores between the suggested and forced selection albeit with more variability 190 across the different scenarios (**Figure** 2c). For example, MutSpec had rF₁ \approx 1 in both medium and 191 hard scenarios indicating that MutSpec is performing worse than SigProfilerExtractor (Figure 192 2a) not because of its algorithm for selecting the total number of signatures but likely due to its 193 implementation of the utilized numerical factorization. SigneR (hard scenarios), 194 SigProfiler PCAWG (hard), SigMiner (medium), TensorSignatures (all), and SigFit (all) had 195 lower F_1 scores for automatic solutions compared to forced solutions ($rF_1 < 1$), thus, indicating 196 that their automatic approaches for selecting the total number of signatures are not optimally 197 performing (Figure 2c). Surprisingly, EMu had higher F₁ scores for automatic solutions in some 198 hard scenarios. Considering the overall performance of EMu (Figure 2a), this outcome likely 199 reflects the lack of convergence during the minimization of the EMu objective function for 200 certain number of signatures in the hard scenarios. 201 202 Overall, across all suggested extractions from noiseless hard scenarios reflecting $\sim 60\%$ of human 203 cancer types, SigProfilerExtractor outperformed all other tools. SigProfilerExtractor was able to 204 identify between 10% and 37% more true positive signatures while yielding between 2.7- and 205 16-fold less false positive signatures compared to the next seven best performing tools: 206 SigProfiler PCAWG, SignatureAnalyzer, SigneR, MutationalPatterns, MutSpec, 207 SomaticSignatures, and SignatureTools (Figure 2*d* and Supplementary Table 1). 208 209

210 Extended benchmarking of SigProfilerExtractor and the other seven top performing tools

211 The reported comparisons for SBS-96 scenarios rely on a cosine similarity ≥0.90 for determining

212 TP signatures and <0.90 for determining FP signatures. Note that a cosine similarity ≥ 0.90 is

highly unlikely to happen purely by chance (p-value = 5.90×10^{-9}) as two random nonnegative

214 vectors are expected to have an average cosine similarity of 0.75 purely by chance³³.

215 Importantly, SigProfilerExtractor's performance does not depend on the specific value of the

cosine similarity threshold (Figure 3a) as the tool consistently outperforms other bioinformatics

approaches for almost any value of the threshold above 0.80 (p-value: 0.057). Cosine similarity

thresholds below 0.80 were not explored as extracted signature may be similar to ground-truth

219 signatures purely by chance.

220

221 Additional benchmarking was performed by generating 12 scenarios simulated using between 3 222 and 30 signatures with an extended number of mutational channels (Supplementary Note 1). 223 SigProfilerExtractor and SignatureAnalyzer are the only two tools that support analysis of 224 custom size matrices and provide GPU support (Table 1), thus, allowing analysis of data with 225 extended number of mutational channels within a reasonable timeframe. In contrast, all other 226 matrix factorization tools rely solely on CPU implementations with full runs expected to take 227 many months for each tool applied to these scenarios (Table 1). SigProfilerExtractor and 228 SignatureAnalyzer exhibited similar performance on the extended noiseless scenarios to that 229 observed on SBS-96 noiseless scenarios. Overall, SigProfilerExtractor outperformed 230 SignatureAnalyzer with average F1 scores of 0.92 and 0.85, respectively (Supplementary Table 231 2).

232

233 To further compare SigProfilerExtractor with the other seven top performing tools, we applied 234 each tool to a dataset with 30 ground-truth SBS-96 signatures operative in 1,000 genomes and 235 random noise between 0% and 10%. Analysis for each noise level was repeated 20 times to 236 account for any variability in the noise generation. SigProfilerExtractor, SomaticSignatures, 237 MutSpec, SignatureToolsLib were robust to noise with mostly unaffected performance (Figure 238 3b and Supplementary Table 3). In contrast, SigProfiler PCAWG, SignatureAnalyzer, SigneR, 239 and MutationalPatterns were susceptible to noise (Figure 3b). For example, 2.5% noise reduced 240 SignatureAnalyzer's F₁ from 0.76 to 0.66 while 10% noise reduced its F₁ to 0.07. Similarly, 10% 241 noise reduced the F_1 of SigProfiler PCAWG from 0.76 to 0.57, the F_1 of SigneR from 0.61 to 242 0.43, and the F₁ of MutationalPatterns from 0.60 to 0.37. SignatureAnalyzer's reduced 243 performance on data with noise is due to its automated approach for selecting total number of 244 signatures. SignatureAnalyzer uses automatic relevance determination³⁴ for selecting the number 245 of signatures with this number increasing from 26 (no noise; 30 ground-truth signatures) to 96 246 signatures (10% noise; Supplementary Table 3). In contrast, SigProfiler PCAWG, SigneR and 247 MutationalPatterns exhibit similar performance between forced and suggested solutions on data 248 with noise (Supplementary Table 3) indicating that their reduced performance is likely due to 249 the numerical implementation of their respective factorization approaches.

250

SigProfilerExtractor outperformed all other tools regardless of the levels of noise. Simulations
with 5% noise reflect genomics datasets with ~0.95 average sensitivity and precision of single
base substitutions, similar to the recently published PCAWG cohort which has 95% sensitivity
(90% confidence interval, 88–98%) and 95% precision (90% confidence interval, 71–99%)²⁸.
For simulations with 5% noise, SigProfilerExtractor was able to identify between 20% and 50%

256	more true positive signatures while yielding more than 5-fold less false positive signatures
257	compared to the next seven best performing tools: SigProfiler_PCAWG, SignatureAnalyzer,
258	SigneR, MutationalPatterns, MutSpec, SomaticSignatures, and SignatureTools (Figure 3c and
259	Supplementary Table 3).
260	
261	Analysis of 2,778 whole-genome sequenced human cancers with SigProfilerExtractor
262	To demonstrate its ability to yield novel biological results, SigProfilerExtractor was applied to
263	the recently published set of 2,778 whole-genome sequenced cancers ²⁸ . As previously done in
264	our original PCAWG analysis of mutational signatures ¹⁴ , extraction of mutational signatures was
265	performed within each cancer type as well as across all samples (Supplementary Data). In
266	addition to all previously detected signatures ¹⁴ , our direct application of SigProfilerExtractor
267	revealed three novel mutational signatures were identified in the PCAWG dataset: SBS92,
268	SBS93, and SBS94 (Figure 4 and Supplementary Table 4).
269	
270	Signature SBS92 was found predominately in PCAWG bladder cancers; the signature was
271	characterized by T>C mutations with strong transcriptional strand bias consistent with damage
272	on purines for all types of single base substitutions (Figure 4a). Signature SBS92 was 9-fold
273	elevated (Figure 4 <i>d</i> ; p-value: 7.6×10^{-3} using Wilcoxon rank sum test) in bladder cancers of ever
274	smokers compared to never-smokers in the PCAWG cohort. An almost identical signature was
275	identified by re-analyzing a recently published cohort of 88 whole-genome sequenced
276	microbiopsies of histologically normal urothelium ³⁵ with the similarity extending to both
277	trinucleotide context and transcriptional strand bias (Figure 4 <i>a</i> ; cosine similarity: 0.98; p-value <

278 10⁻²⁵⁶). Consistently, SBS92 was found to be 3-fold elevated in the normal urothelium of tobacco

ever smokers compared to never-smokers (Figure 4*d*; p-value: 8.3 x 10⁻³ using Wilcoxon rank
sum test).

282	Signature SBS93 was identified almost exclusively in PCAWG stomach cancers. SBS93 was
283	characterized by T>C and T>G mutations with a strand bias consistent with damage on
284	pyrimidines for TpTpA contexts (mutated base underlined; Figure 4b). De novo extraction from
285	the Mutographs cohort of 552 whole-genome sequenced esophageal squamous cell carcinomas ³⁶ ,
286	a cancer type not included in the PCAWG dataset ²⁸ , identified an analogous mutational signature
287	with the similarity extending to both trinucleotide context and transcriptional strand bias (Figure
288	4b; cosine similarity: 0.88; p-value: 1.1 x 10 ⁻⁶). Signature SBS94 was found at high levels in a
289	single colorectal PCAWG cancer with smaller contributions to another 8 colorectal cancers. The
290	pattern of SBS94 was characterized by C>A mutations with a strand bias indicative of damage
291	on guanine (Figure $4c$). Validation of somatic mutations by visual inspection confirmed that
292	98% of mutations contributed by SBS94 are likely real. Signatures SBS93 and SBS94 did not
293	associate with any of the available PCAWG metadata ²⁸ and their etiologies remain unknown.
294	

295 **DISCUSSION**

296 The performed large-scale benchmarking demonstrates that SigProfilerExtractor outperforms 297 thirteen other tools for *de novo* extraction of mutational signatures for noiseless datasets as well 298 as for datasets containing matrices with different levels of random noise. Importantly, 299 SigProfilerExtractor generates almost no false positive signatures while still identifying a higher 300 number of true positive signatures when compared to any of the other tools (Figure 2d and 301 Figure 3c). De novo extraction of mutational signatures relies both on a factorization approach 302 and on a model selection algorithm for determining the total number of operative signatures 303 (Figure 1). Benchmarking with forced model selection, where tools were required to extract the 304 known number of ground-truth mutational signatures, reveals that SigProfilerExtractor's 305 factorization performs better when compared to the factorizations of other tools (Figure 2a and 306 Supplementary Tables 1-3). Similarly, benchmarking with suggested model selection, which 307 most closely matches analysis of a real dataset with unknown number of signatures, further 308 demonstrates SigProfilerExtractor's ability to reveal novel biological results (Figure 2a and 309 Supplementary Tables 1-3).

310

311 While our benchmarking evaluated thirteen additional tools, six of the thirteen tools internally 312 rely on the same computational engine. Maftools, MutationalPatterns, MutSpec,

SignatureToolsLib, SigMiner, and SomaticSignatures use the NMF R package³⁷ to perform their factorization (**Table 1**), albeit with slightly different hyperparameters and, in some cases, distinct pre-processing of the input matrix. Predictably, these six tools have similar performance across many of the scenarios (**Supplementary Tables 1-3**). SigProfiler_PCAWG and MutSignatures utilize customize versions of an NMF implementation originally developed by Brunet *et al.*³⁸ for

318	analysis of gene expression data. TensorSignatures makes use of the standard factorization
319	algorithms included in TensorFlow ³⁹ . SigFit uses a previously developed nonnegative
320	factorization method, viz., Stan R package ⁴⁰ . In contrast, EMu, SignatureAnalyzer, SigneR, and
321	SigProfilerExtractor provide original implementations of their factorization algorithms (Table
322	1). EMu was originally developed and tested on small datasets (<i>e.g.</i> , 21 breast genomes) ¹⁵ and its
323	benchmarking performance is perhaps unsurprising considering the large number of synthetic
324	samples used in all scenarios. Surprisingly, the original implementations of SignatureAnalyzer
325	and SigneR were susceptible to noise, yielding high numbers of false-positive signatures (Figure
326	3 <i>b</i>).
327	
328	Seven of the tools did not provide an automatic approach for selecting the total number of
329	operative signatures in a dataset (Table 1). Instead, most of these tools offered methodologies for
330	manually selecting the optimal number of signatures bringing user-dependence and arbitrariness
331	in selecting solutions. EMu, SigFit, SigMiner, SignatureAnalyzer, SigneR, TensorSignatures,
332	and SigProfilerExtractor provided capabilities for automatically selecting the total number of
333	operative signatures. EMu, TensorSignatures, SigneR select the total number of signatures using
334	Bayesian information criterion (BIC) ⁴¹ , while SignatureAnalyzer and SigMiner utilize automatic
335	relevance determination (ARD) ³⁴ . SigFit's selection approach is based on the Elbow method ⁴² .
336	SigProfilerExtractor leverages a modified version of the NMFk selection approach which was
337	previously tested on more than 55,000 synthetic random matrices with pre-determined latent
338	factors and shown to outperform other model selection approaches ⁴³ . Importantly, our
339	simulations demonstrate that SigProfilerExtractor's model selection is robust to noise while the

implemented BIC and ARD approaches are affected even by low levels of noise (Figure 3b).

341 In addition to outperforming thirteen other tools on simulated datasets, SigProfilerExtractor can 342 reveal additional biological results as demonstrated by identifying three novel signatures from 343 reanalysis of the PCAWG dataset. Importantly, SigProfilerExtractor identifies signature SBS92 344 (Figure 4) which is associated with tobacco smoking in whole-genome sequenced bladder 345 cancers and in whole-genome sequenced microbiopsies from normal bladder urothelium. The 346 strong transcriptional strand bias observed in SBS92 is indicative of an environmental mutagen 347 exposure that damages purines. Tobacco smoke is a complex mixture of at least 60 chemicals¹³, 348 many capable of causing damage on purines. Interestingly, our and other prior analyses of exome sequenced bladder cancers from The Cancer Genome Atlas (TCGA) project^{13,44} did not reveal 349 350 SBS92. Reanalysis of the set of TCGA bladder cancer exomes⁴⁵ with SigProfilerExtractor was 351 also unable to detect SBS92 (Supplementary Data). We suspect that the lack of SBS92 in the 352 TCGA bladder cancers was due to the use of exome sequencing; note that SBS92 is 353 predominately found in intergenic regions (Figure 4a) with most samples expected to have less 354 than 15 mutations from SBS92 in their exomes. To confirm this hypothesis, we downsampled the 355 whole-genome sequenced bladder cancers and the whole-genome sequenced microbiopsies from 356 normal bladder urothelium to exomes. SigProfilerExtractor's analysis of these downsampled 357 genomes was unable to detect SBS92 confirming that exome sequencing is insufficient to 358 identify signature SBS92 (Supplementary Data).

359

360 In summary, here we report SigProfilerExtractor – a computational tool for *de novo* extraction of 361 mutational signatures. We demonstrate that SigProfilerExtractor outperforms thirteen other tools 362 by conducting the largest benchmarking of bioinformatics approaches for extracting mutational 363 signatures. Further, we apply SigProfilerExtractor to 2,778 whole-genome sequenced cancers

- 364 and reveal several novel mutational signatures including a signature putatively attributed to
- 365 tobacco smoking mutagenesis in bladder cancer and in normal bladder epithelium.

367 ONLINE METHODS

368 **Computational implementation of SigProfilerExtractor and its seven modules**

369 The implementation of SigProfilerExtractor can be separated into seven distinct modules which 370 are packaged together into a single bioinformatics tool. *Module 1* processes the initial input data, 371 which can be provided as either a mutational catalogue containing a set of somatic mutations or a 372 mutational matrix. Module 2 is responsible for resampling and normalization of the mutational 373 matrix prior to performing nonnegative matrix factorization. Module 3 performs matrix 374 factorization using nonnegative matrix factorization with multiple replicates. Module 4 utilizes 375 custom clustering to derive consensus solutions and to perform model selection. Module 5 376 decomposes the derived set of *de novo* signatures to a set of previously derived COSMIC 377 signatures. Module 6 is responsible for calculating the activities of different signatures in 378 individual samples. *Module* 7 handles the extensive outputting and plotting of the different 379 analysis performed by SigProfilerExtractor. In principle, each of these modules allows extensive 380 customization. SigProfilerExtractor provides a seamless integration of these seven modules that 381 allows using them in an orchestrated and preconfigured manner with little input from a user.

382

383 *Module 1: Processing of input mutational catalogues or input mutational matrices*

SigProfilerExtractor deciphers mutational signatures from a mutational matrix M with t rows and n columns; rows represent mutational channels while columns reflect individual cancer samples (**Figure 1**a). The value of each cell in the matrix, M, corresponds to the number of somatic mutations from a particular mutational channel in each sample. The mutational matrix can be provided as a text file with the first column containing the names of the mutational channels and the first row containing the names of the examined samples. Alternatively, users

can provide a mutational catalogue of somatic mutations in a commonly used format (*e.g.*, VCF,
 MAF, *etc.*) and this mutational catalogue will be internally converted into the appropriate
 mutational matrix by SigProfilerMatrixGenerator².

393

394 Module 2: Resampling of the input mutational matrix and normalizing the resampled matrix

395 SigProfilerExtractor does not factorize the original input matrix. Rather, prior to performing 396 matrix factorization, SigProfilerExtractor performs independent Poisson resampling of the 397 original matrix for each replicate¹. As such, the matrix factorized in each replicate is never the 398 same for a given value of k (Figure 1b). The resampling is performed to ensure that Poisson 399 fluctuations of the matrix do not impact the stability of the factorization results. Additional 400 normalization is performed after resampling to overcome potential skewing of the factorization 401 from any hypermutators. SigProfilerExtractor supports four standard normalization methods⁴⁶: 402 (i) Gaussian mixture model (GMM) normalization (default); (ii) 100X normalization; (iii) log2 403 normalization; (iv) no normalization. No normalization does not perform any additional 404 transformation on the Poisson resampled matrix. In log2 normalization, the sum of each column 405 in the matrix is derived and logarithm with base 2 is calculated for each of these sums. Each cell 406 in a column of the matrix is multiplied by the log2 of the column-sum and subsequently divided by the original column sum. In 100X normalization, the sum of each column in the matrix is 407 408 derived. For each column where the sum exceeds 100 times the number of mutational channels 409 (*i.e.*, 100 times the number of rows in the matrix), each cell in the column is multiplied by the 410 100 times the number of mutational channels and subsequently divided by the original column 411 sum. This normalization ensures that no sample has a total number of mutations above 100 times 412 the number of mutational channels. GMM normalization encompasses a two-step process. The

413 first step derives the normalization cutoff value in a data-driven manner using a Gaussian 414 mixture model (GMM). The second step normalizes the appropriate columns using the derived 415 cutoff value. The first step uses a GMM to separate the samples into two groups based on their 416 total number of mutations; the total number of mutations in a sample reflects the sum of a 417 column in the matrix. The group with larger number of samples is subsequently selected, and the 418 same process is applied iteratively until it converges. Convergence is achieved when the mean of 419 the two groups is separated by no more than four standard deviations of the larger group. A 420 cutoff value is derived as the average value plus two standard deviations from the total number 421 of somatic mutations in the last large group. If the derived cutoff value is below 100 times the 422 number of mutational channels, the cutoff value is adjusted to 100 times the number of mutational channels. For each column where the sum exceeds the derived cutoff value, each cell 423 424 in the column is multiplied by the cutoff value and subsequently divided by the original column 425 sum. Note that no normalization is performed if the means of the first two groups are not 426 separate by at least four standard deviations. In all cases, columns with a sum of zero, reflecting, 427 genomes without any somatic mutations, are ignored to avoid division by zero.

428

429

430 Module 3: Matrix Factorization Using Nonnegative Matrix Factorization with Replicates

By default, SigProfilerExtractor factorizes the matrix M with different ranks searching for an optimal solution between k=1 and k=40 mutational signatures. For each value of k, by default, the tool performs 500 independent nonnegative matrix factorizations of the normalized Poisson resampled input matrix. Thus, for each value of k, SigProfilerExtractor generates 500 distinct factorizations of normalized Poisson resampled matrices resulting into 500 different matrices S, 436 each matrix reflecting the patterns of the *de novo* mutational signatures, and 500 different 437 matrices A, each matrix reflecting the activities of the *de novo* mutational signatures (Figure 1b). 438 To perform each of these factorizations, SigProfilerExtractor utilizes a custom implementation of 439 the multiplicative update algorithm³. Specifically, SigProfilerExtractor initializes the S and A440 matrices in the first step of the factorization using either random initial conditions (default) or 441 one of the derivatives of nonnegative double singular vector decomposition⁴⁷. 442 SigProfilerExtractor provides internal support for minimizing three different objective functions 443 based on: (i) generalized Kullback-Leibler updates (default); (ii) Euclidean updates; (iii) Itakura-444 Saito updates. By default, the tool performs all factorization using multithreading of central 445 processing units (CPUs) and provides support for factorization using graphics processing units (GPUs) by leveraging PyTorch⁴⁸. In all cases, by default, the implemented minimization 446 447 performs at least 10,000 iterations (also known as NMF updates or NMF multiplicative update 448 steps) with a maximum of 1,000,000 iterations. By default, the convergence tolerance of the algorithm is set to 10⁻¹⁵. Note that SigProfilerExtractor allows configuring all factorization 449 450 parameters.

451

452 Module 4: Custom partition clustering and performing model selection

The previously described *Module 3* generates a number of sets with each set containing, by default, 500 different matrices S, where each matrix reflects the patterns of *de novo* mutational signatures for a particular factorization of a normalized Poisson resampled matrix. One set, containing 500 different matrices S, is generated for each of the interrogated total number of operative signatures, k, with a default range for k between 1 and 40 signatures. For each value of k, *Module 4* first performs custom clustering of the S matrices and, subsequently, applies a

459 modified version of the NMFk model selection approach to select the optimal value of k^{43} 460 (Figure 1b). Specifically, for each value of k, the clustering is initialized with k random 461 centroids. One of the S matrices is randomly chosen, and its columns matched to the most similar 462 centroids with no two columns assigned to the same cluster. The process is repeated until the 463 columns of all S matrices in the set are assigned to their respective clusters. SigProfilerExtractor implements the Hungarian algorithm³⁰ to pair consensus vectors from two matrices (*i.e.*, cluster 464 465 centroids and mutational signature from a matrix \boldsymbol{S} ; the Hungarian algorithm maximizes the 466 total cosine similarities of all paired vectors between two matrices³⁰. After assigning all columns 467 to a cluster, the centroids of each cluster are recalculated by evaluating the average of all 468 columns/vectors in a cluster. This process continues iteratively until the average silhouette coefficient converges (*i.e.*, its value does not change by more than 10^{-12}). After convergence for a 469 470 given value of k, the centroids of the clusters are reported as consensus mutational signatures, an 471 overall reconstruction error is calculated for describing the original input matrix, M, and stability is calculated for each signature by computing the silhouette value⁴⁹ of the cluster corresponding 472 473 to that signature (Figure 1b). The silhouette value of a cluster measures the similarities of the 474 objects assigned to that cluster compared to any other cluster. Silhouette values range from -1.0 475 to +1.0 with values above zero indicating that, on average, objects have a higher similarity with 476 their own cluster compared to their nearest clusters. Note that signatures with low stability 477 correspond to a lack of uniqueness of the NMF due to Poisson resampling and/or to the potential 478 existence of multiple convergent stationary points in the NMF solution³².

479

480 Our custom clustering is performed for each of the interrogated total number of operative

481 signatures, *k*, with a default range for *k* between 1 and 40 signatures. After performing clustering,

482 for each value of k, one has derived: (i) the consensus set of mutational signatures; (ii) an overall 483 reconstruction error for describing the original input matrix; and (iii) stability value for each of 484 the identified consensus mutational signatures.

485

486 SigProfilerExtractor performs a solution selection based on the stability of signatures in a 487 solution and the ability of these signatures to reconstruct the original input matrix. By default, 488 SigProfilerExtractor will consider solutions stable if the signatures derived in the solution have 489 an average stability above 0.80 with no individual signature having stability below 0.20. To 490 reduce overfitting, the tool also measures the information gained from the extracted set of 491 signatures in each solution. SigProfilerExtractor compares, using Wilcoxon rank-sum tests, the 492 reconstruction errors across all samples from the stable solution with the greatest number of 493 signatures to the reconstruction errors across all samples from stable solutions with lower 494 number of signatures. Stable solutions with lower number of signatures are compared in a 495 decreasing order to their total number of signatures with comparison stopping if the Wilcoxon 496 rank-sum test yields a *p*-value below 0.05 (*i.e.*, reflecting that a solution does not describe the 497 original data as good as the stable solution with the greatest number of signatures). The stable 498 solution with lowest number of signatures and a Wilcoxon rank-sum test *p-value* above 0.05 is 499 selected as the optimal solution. If no solution has a Wilcoxon rank-sum test *p-value* above 0.05, 500 the stable solution with the greatest number of signatures is selected as the optimal solution. Note 501 that while SigProfilerExtractor selects an optimal solution, it outputs all the information 502 necessary to evaluate mutational signatures and their activities for all other stable and unstable 503 solutions.

505 Module 5: Decomposing de novo extracted signatures to known COSMIC signatures

506 SigProfilerExtractor provides a module for decomposing each of the *de novo* extracted 507 mutational signatures to a set of previously derived signatures. By default, the tool decomposes each of the signatures in the optimal solution to a set of COSMICv3 reference signatures¹⁴ with 508 509 support for signatures of single base substitutions (SBS), doublet base substitutions (DBS), and 510 small insertions and deletions (ID). Since the SBS COSMICv3 reference signatures were derived 511 under the SBS-96 classification², any extended classification of single base substitutions (e.g., SBS-288 and SBS-1536)² is first collapsed to the SBS-96 classification and, subsequently, 512 513 decomposed to the COSMICv3 reference signatures¹⁴. The decomposition functionality leverages nonnegative least square (NNLS) algorithm⁵⁰ as its main computational engine. A 514 515 mixture of addition and removal steps (add-remove functionality) were developed to estimate the 516 list of COSMIC signatures for a *de novo* signature. Specifically, for each *de novo* signature, a 517 COSMIC signature is iteratively added to a list of signatures used to explain the *de novo* 518 signature, NNLS is applied, and the signature which addition causes the greatest decrease of the 519 L2 error is selected. If this greatest decrease is more than a specific threshold (default value of 520 0.05) then the signature is included in the list of signatures used to explain the *de novo* signature. 521 The addition is immediately followed by a removal step. Each COSMIC signature in the list of 522 signatures used to explain the *de novo* signature are iteratively removed, NNLS is applied, and 523 the signature that causes the least decrease of the L2 error is selected. If this least decrease is less 524 than a specific threshold (default value of 1%) then the signature is removed from the list of 525 signatures used to explain the *de novo* signature. The addition and removal steps are iterated until 526 no signatures are added or removed from the list of signatures used to explain the *de novo* 527 signature. Several previously implemented rules for mutational signatures are incorporated by

528	default in the decomposition module ¹⁴ . Specifically, for signatures of single base substitutions:
529	(i) the list of signatures used to explain the <i>de novo</i> signature is initialized with clock-like
530	signatures SBS1 and SBS5; ¹¹ (ii) biologically connected signatures are included as previously
531	done in Ref ¹⁴ (<i>e.g.</i> , if SBS17a is included in the list then SBS17b is also included the list). The
532	decomposition module is highly customizable as it allows changing all default parameters as
533	well as adding additional new rules or removing existing rules for inclusion and exclusion of
534	particular signatures.

535

536 Module 6: Evaluating activities of mutational signatures in individual samples

537 De novo extracted and COSMIC derived signatures are refitted to individual samples using nonnegative least squares (NNLS)⁵⁰. *Module 6* internally utilizes the add-remove functionality of 538 539 *Module 5* with each sample in the original matrix, *M*, being individually examined. For *de* 540 novo mutational signatures, all *de novo* signatures are initially added to the list of signatures used 541 to explain the sample and a removal step with a cutoff of 2% is applied. To assign COSMIC 542 signatures in a sample, the module first derives the set of *de novo* signatures in that sample. 543 Decomposition to the COSMICv3 signatures using *Module 5* is performed for each of the *de* 544 novo signatures and the identified COSMICv3 signatures are refitted using the add-remove 545 functionality with a removal and addition cutoffs set at 5%. Finally, the activity matrix is 546 constructed by combining the activity vectors generated for all samples in the dataset.

547

548 Module 7: Outputting and plotting of analysis results

549 All previous modules make use of *Module* 7 for outputting and plotting of the generated results.

550 It should be noted that SigProfilerExtractor provides extensive output for the interrogated total

551	number of operative signatures, k , with a default range of k between 1 and 40 signatures. For
552	each value of k, SigProfilerExtractor outputs the set of operative de novo mutational signatures,
553	the activities of the operative signatures, and an extensive set of information related to individual
554	samples, individual de novo signatures, and the overall convergence of the factorization and
555	clustering. Module 7 also provides additional information when ran in debug mode. In addition
556	to outputting information, SigProfilerExtractor also generates a bouquet of plots both for each
557	value of k as well as for the suggested optimal solution. SigProfilerExtractor utilizes all
558	previously implemented plots in SigProfilerPlotting ² as well as includes several newly developed
559	plots.
560 561 562	Analysis of the genomics data from cancer and normal somatic tissues
563	For all examined cancer and normal somatic tissues, de novo extraction of mutational signatures
564	was performed with SigProfilerExtractor with default parameters using two distinct mutational
565	classifications: SBS-96 and SBS-288. The SBS-96 mutation classification incorporates the six
566	types of single base substitutions: C>A, C>G, C>T, T>A, T>C, and T>G. Each type of single
567	base substitution is further separated into 16 subtypes determined by the four possible bases 5'
568	and 3' adjacent to each mutated base. The SBS-288 mutation classification extends the SBS-96
569	mutation classification by adding additional information for each of the 96 subtypes.
570	Specifically, SBS-288 incorporates whether a single base substitution is in non-
571	transcribed/intergenic DNA, on the transcribed strand of a gene, or on the untranscribed strand of
572	the gene. De novo extraction was performed separately for all examined datasets. Specifically,
573	SigProfilerExtractor was applied: (i) to all 2,778 whole-genome sequenced cancers from the Pan-
574	Cancer Analysis of Whole Genomes project ²⁸ ; (ii) to all samples in each of the 37 cancer types

575 of Pan-Cancer Analysis of Whole Genomes project²⁸ with each cancer type examined separately; 576 *(iii)* to all 88 whole-genome sequenced microbiopsies of histologically normal urothelium³⁵; *(iv)* to the complete set of bladder cancers from TCGA⁴⁵; (v) to exome downsampling of all bladder 577 578 whole-genome sequenced cancers from the Pan-Cancer Analysis of Whole Genomes project²⁸; 579 (vi) to exome downsampling of all 88 whole-genome sequenced microbiopsies of histologically 580 normal urothelium³⁵. In all cases, the mutational catalogue of each sample was taken from the 581 respective original publications. The results from all performed *de novo* extractions can be found 582 in Supplementary Data. Downsampling of a whole-genome sequenced sample to a whole-583 exome was performed using SigProfilerMatrixGenerator². 584 585 Additional approaches for miscellaneous analysis 586 Synthetic scenarios were labeled as easy, medium, and hard based on the number of operative 587 signatures in each scenario. Based on our most recent analysis of mutational signatures in 82 588 cancer types¹⁴, approximately 7.4% of human cancer types have 5 or less signatures (reflected in 589 simulations of easy scenarios), 15.9% have 11 to 21 signatures (medium scenarios), and 59.5% 590 have 25 or more signatures (hard scenarios). Note that 17.2% of cancer types have either 591 between 5 and 10 signatures or between 22 and 24 signatures. 592 593 Cosine similarity was used to compare the profiles of different mutational signatures. P-values 594 can be attributed to cosine similarities based on a null hypothesis of uniform random distribution 595 of nonnegative vectors³³.

597	Briefly, the prevalence of somatic mutations in a whole-exome sample was calculated based on
598	the identified mutations in protein coding genes and assuming that an average whole-exome has
599	sufficient coverage of 30.0 megabase-pairs in protein coding genes. The prevalence of somatic
600	mutations in a whole-genome sample was calculated based on all identified mutations and
601	assuming that an average whole-genome has sufficient coverage of 3.00 gigabase-pairs.
602	
603	All methods related to the generation of the benchmarking scenarios and the application of the
604	different tools to these scenarios can be found in Supplementary Note 1.

606

607 TABLES

			Factorization Approach			Selection Approach			Supported Contexts		COSMIC
Tool Name	Input	Platform	Method	Computational Engine	GPU	Manual	Automatic	Automatic Algorithm	Mutational Catalogue Support	Plotting	Comparison
EMu ¹⁵	Matrix	C++	EM	Original implementation ¹⁵	No	Yes	Yes^	BIC ⁴¹	SBS-96	No	No
Maftools ¹⁶	Matrix MAF	R-Bioconductor	NMF	NMF R package ³⁷	No	Yes	No	-	SBS-96	SBS-96	1-to-1
MutationalPatterns ¹⁷	Matrix VCF	R-Bioconductor	NMF	NMF R package ³⁷	No	Yes	No	-	SBS-96 SBS-192	SBS-96 SBS-192	1-to-1
MutSignatures ¹⁸	VCF MAF Matrix	R	NMF	Brunet et al. ³⁸	No	No	No	-	SBS-96	SBS-96	1-to-1
MutSpec ¹⁹	Matrix VCF Custom	Galaxy Perl R	NMF	NMF R package ³⁷	No	Yes	No	-	SBS-96 SBS-192	SBS-96 SBS-192	1-to-1
SigFit ²⁰	Matrix	R	Bayesian Inference	Stan R package ⁴⁰	No	Yes	Yes^	Elbow method ⁴²	SBS-96	SBS-96 SBS-192	1-to-1
SigMiner ²¹	Matrix MAF	R	[automatic] Bayesian NMF [manual] NMF	[automatic] SignatureAnalyzer implementation ²² [manual] NMF R package ³⁷	No	Yes^	Yes	ARD ³⁴	SBS-96 DBS-78 ID-83	Generic	1-to-1
SignatureAnalyzer ^{22,23}	Matrix MAF	R [CPU] ¹⁸ Python [GPU] ¹⁹	Bayesian NMF	Original implementation ^{22,23}	Yes	No	Yes	ARD ³⁴	SBS-96 DBS-78 ID-83	SBS-96 DBS-78 ID-83	1-to-1
SignatureToolsLib ²⁴	Matrix VCF Custom	R	NMF	NMF R package ³⁷	No	Yes	No	-	SBS-96 DBS-78 ID-83 SV-32	SBS-96 SV-32 Generic	No
SigneR ²⁵	Matrix VCF	R-Bioconductor C++	Bayesian NMF	Original implementation ²⁵	No	Yes	Yes^	BIC ⁴¹	SBS-96	SBS-96	No
SigProfilerExtractor	Matrix VCF MAF Custom	Python R wrapper	NMF	[current report] Original implementation	Yes	Yes	Yes^	NMFk ⁴³	SBS-96 DBS-78 ID-83 CN-48 Others ² Any	SBS-96 DBS-78 ID-83 CN-48 SV-32 Others ² Generic	1-to-many
SigProfiler_PCAWG ¹⁴	Matrix VCF MAF Custom	Python MATLAB	NMF	Brunet et al. ³⁸	No	Yes	No	-	SBS-96 DBS-78 ID-83 Others ² Any	SBS-96 DBS-78 ID-83	No
SomaticSignatures ²⁶	Matrix VCF	R-Bioconductor	NMF PCA	NMF R package ³⁷ pcaMethods R package ⁵¹	No	Yes	No	-	SBS-96	SBS-96	No
TensorSignatures ²⁷	VCF	Python	NTF	TensorFlow ³⁹	Yes	Yes	Yes^	BIC ⁴¹	Tensor	SBS-96 with strand bias	No

609	Table 1: Overview of bioinformatics tools for <i>de novo</i> extraction of mutational signatures.
610	Tools are ordered alphabetically. Notations: ^ denotes the default approach for selecting the total
611	number of signatures when a tool supports both manual and automatic selection; 1-to-1 refers to
612	one de novo signature being matched with exactly one COSMIC signature; 1-to-many refers to
613	one <i>de novo</i> signature being matched with a combination of one or more COSMIC signatures.
614	Abbreviations: MAF: mutation annotation format; VCF: variant call format; EM: expectation-
615	maximization algorithm; NMF: nonnegative matrix factorization; NTF: nonnegative tensor
616	factorization; ARD: automatic relevance determination; BIC: Bayesian information criterion;
617	COSMIC: catalogue of somatic mutations in cancer; SBS: single base substitutions; DBS:
618	doublet base substitutions; ID: small insertions and deletions; SV: structural variants.
619	

620 FIGURE LEGENDS

621 Figure 1. Overview of SigProfilerExtractor. (a) SigProfilerExtractor's general workflow is 622 outlined starting from an input of somatic mutations and resulting in an output of *de novo* 623 mutational signatures. An example is shown for a solution with three *de novo* signatures. 624 Somatic mutations are first converted into a mutational matrix. Subsequently, the matrix is 625 factorized with different ranks using nonnegative matrix factorization. Model selection is applied 626 to identify the optimal factorization rank based on each solution's stability and its reconstruction 627 of the original data. (b) Schematic representation for an example decomposition with a 628 factorization rank of k=3 reflecting three operative mutational signatures. By default, 629 SigProfilerExtractor performs 500 independent nonnegative matrix factorizations with the matrix 630 M being Poisson resampled and normalized (denoted by ^) prior to each factorization. Partition 631 clustering of the 500 factorizations is used to evaluate the factorization stability rank, measured 632 in silhouette values; clustering can also be presented as two-dimensional projections revealing 633 more similar mutational signatures as shown for the three example signatures. The centroid of 634 the clustered solutions (denoted by -) is compared to the original matrix *M*. (c) All identified de 635 novo signatures are matched to a combination of known COSMIC mutational signatures. An 636 example is given for *de novo* extracted signature SBS96B which matches a combination of 637 COSMIC signatures SBS1, SBS2, and SBS13.

638

639 Figure 2. Benchmarking of bioinformatics tools for *de novo* extraction of mutational

640 signatures using SBS-96 noiseless scenarios. (a) Average precision (x-axes), sensitivities (y-

641 axes), and F₁ scores (harmonic mean of precision and sensitivity; red curves) are shown across

642 the three types of scenarios. Different tools are displayed using circles and triangles with

643 different colors. Circles are used to display results for suggested model selection, which most 644 closely matches analysis of a real dataset. Triangles are used to display results for forced model 645 selection, where tools were required to extract the known total number of ground-truth 646 mutational signatures. All triangles are located on the diagonal as the forced model selection 647 results in equal numbers of false positive and false negative signatures. (b) Evaluating the effect 648 of ground-truth signatures on the *de novo* extraction by different tools (x-axes). Ratio of F₁ 649 scores (y-axes) with confidence intervals were calculated for medium complexity scenarios 650 simulated using COSMIC, SA, or random signatures. Ratio of approximately 1.00 indicates a 651 similar performance between different types of signatures. (c) Evaluating the performance of de 652 *novo* extraction between suggested and forced selection for different tools (x-axes). Ratio of F_1 653 scores (y-axes) with confidence intervals were calculated for all medium and hard scenarios. 654 Ratio of approximately 1.00 indicates a similar performance between suggested and forced 655 model selection. (d) Summary of the performance for the top seven tools on hard SBS-96 656 noiseless scenarios with suggested model selection. Y-axes reflect F_1 score (left plot), sensitivity 657 (middle plot), and false discovery rate (right plot), respectively. Results from SignatureAnalyzer 658 and MutSignatures are not displayed in panels (a), (b), and (c) for forced and suggested model 659 selections, respectively, as the tools do not support these types of analyses.

660

Figure 3. Additional evaluations of the top seven bioinformatics tools for *de novo* extraction of mutational signatures. (*a*) Average F₁ scores for the top seven tools based on different thresholds for cosine similarity in suggested medium and hard scenarios; thresholds for cosine similarity are used for determining true positive signatures (Supplementary Figure 1). X-axes reflect the cosine similarity thresholds, while the Y-axes correspond to the average F₁ scores

666	corresponding to cosine similarity thresholds. (b) Precision and sensitivity of the top seven tools
667	for SBS-96 scenarios with different levels of noise. Noise levels reflect the average number of
668	somatic mutations in a cancer genome affected by additive white Gaussian noise; for example,
669	1% noise corresponds to approximately 1% of mutations in a sample being due to noise. (c)
670	Summary of the performance of the top seven tools on SBS-96 scenarios with 5% noise. Y-axes
671	reflect F1 score (left plot), sensitivity (middle plot), and false discovery rate (right plot),
672	respectively.

673

682

674 Figure 4. Novel signatures identified in the PCAWG cohort of 2,778 whole-genome

675 **sequenced cancers.** Mutational signatures are displayed using 96-plots. Single base substitutions

are shown using the six subtypes of substitutions: C>A, C>G, C>T, T>A, T>C, and T>G.

677 Underneath each subtype are 16 bars reflecting the sequence contexts determined by the four

678 possible bases 5' and 3' to each mutated base. Additional information whether mutations from a

679 signature are in non-transcribed/intergenic DNA, on the transcribed strand of a gene, or on the

680 untranscribed strand of the gene is provided adjacent to the 96 plots. (a) Mutational profile of

681 signature SBS92 derived from the PCAWG cohort (top). Confirmation of the profile of signature

683 histologically normal urothelium³⁵. *(b)* Mutational profile of signature SBS93 derived from the

SBS92 (bottom) by analysis of an independent whole-genome sequenced set of microbiopsies of

684 PCAWG cohort (top). Confirmation of the profile of signature SBS93 (bottom) by analysis of an

684 PCAWG cohort (top). Confirmation of the profile of signature SBS93 (bottom) by analysis of an

685 independent whole-genome sequenced set of esophageal squamous cell carcinomas²⁸. (c)

686 Mutational profile of signature SBS94 derived from the PCAWG cohort. Signature SBS94 was

687 not identified in any additional independent cohort. (d) Bars are used to display average values

688 for numbers of somatic substitutions per megabase (Mb) attributed to signature SBS92 in bladder

- 689 cancer and normal bladder urothelium. Green bars represent never-smokers, whereas blue bars
- 690 correspond to ever smokers. Error bars correspond to 95% confidence intervals. Each p-value is
- 691 based on a Wilcoxon rank sum test.

693 SUPPLEMENTARY INFORMATION

694 Supplementary Figure 1. Standard set of performance metrics used for benchmarking all

- 695 bioinformatics tools. An example demonstrating the derivation of *true positive* (TP), *false*
- 696 positive (FP), or false negative (FN) signatures for a tool applied to a synthetic dataset generated
- 697 using 6 ground truth signatures (termed, Ground Truth Signatures 1 through 6). The tool extracts
- 698 4 signatures (termed, Extracted Signatures A through D). In this example, an extracted signature
- 699 is considered a true positive if it matches one of the ground-truth signatures with a cosine
- 700 similarity threshold of at least 0.90.
- 701

702 Supplementary Table 1. Detailed performance metrics after applying each tool across all

703 SBS-96 noiseless synthetic scenarios. Performance metrics are calculated as per *Supplementary*

Figure 1. An extracted signature is considered a true positive if it matches one of the ground-

truth signatures with a cosine similarity threshold of at least 0.90.

- 706
- 707 Supplementary Table 2. Detailed performance metrics after applying the seven best
- 708 performing tools across SBS-96 synthetic scenarios with different levels of noise.

709 Performance metrics are calculated as per *Supplementary Figure 1*. An extracted signature is

- considered a true positive if it matches one of the ground-truth signatures with a cosine similarity
- 711 threshold of at least 0.90.

- 713 Supplementary Table 3. Detailed performance metrics of applying SigProfilerExtractor
- 714 and SignatureAnalyzer to extended synthetic scenarios. Performance metrics are calculated

715	as per Supplementary Figure 1. An extracted signature is considered a true positive if it matches
716	one of the ground-truth signatures with a cosine similarity threshold of at least 0.90.
717	
718	Supplementary Table 4. Profiles of three novel mutational signatures identified in the
719	PCAWG cohort of 2,778 whole-genome sequenced cancers. The profiles of the novel
720	mutational signatures are reported using the SBS-288 classification which incorporates the
721	trinucleotide context and strand information (intergenic region, untranscribed strand, or
722	transcribed strand) for each type of single base substitution. The SBS-288 classification can be
723	easily collapsed to the commonly used SBS-96 classification.
724	
725	Supplementary Note 1. Detailed description of the performed benchmarking. The
726	supplementary note provides extensive details about each of the generated synthetic scenarios as
727	well as about applying each of the tools to these scenarios. The results from applying all tools to
728	all scenarios, including appropriate input and out files, can be found in Supplementary Data.
729	
730	Supplementary Data
731	All results from the benchmarking with synthetic datasets, including the appropriate input used
732	to run each of the tools as well as the output generated by each of the tools, can be found at:
733	ftp://alexandrovlab-ftp.ucsd.edu/pub/publications/Islam_et_al_SigProfilerExtractor/Benchmark/.
734	
735	All results from the <i>de novo</i> extraction of mutational signatures from the PCAWG dataset can be
736	found at: <u>ftp://alexandrovlab-</u>
737	ftp.ucsd.edu/pub/publications/Islam_et_al_SigProfilerExtractor/PCAWG_Reanalysis/.

738

- All results from the *de novo* extraction of mutational signatures for confirming the patterns of the
- 740 novel signatures for additional datasets can be found at: <u>ftp://alexandrovlab-</u>
- 741 ftp.ucsd.edu/pub/publications/Islam et al SigProfilerExtractor/Confirmation of Novel Signatu

742 <u>res/</u>.

- 743
- All results from the *de novo* extraction of mutational signatures from downsampling of whole-
- genome sequenced samples to whole-exomes can be found at: <u>ftp://alexandrovlab-</u>
- 746 <u>ftp.ucsd.edu/pub/publications/Islam et al SigProfilerExtractor/Downsampling of whole geno</u>
- 747 <u>mes/</u>
- 748

749 ACKNOWLEDGEMENTS

750 The authors would like to thank Allan Balmain (UC San Francisco) for the many useful

- 751 discussions as well as Ville Mustonen (University of Helsinki) and Israel Tojal Da Silva (A.C.
- 752 Camargo Cancer Center) for help in configuring EMu and SigneR, respectively. This work was
- supported by Cancer Research UK Grand Challenge Award C98/A24032 (LBA, PB, and MRS),
- 754 Wellcome grant reference 206194 (MRS), as well as US National Institute of Health grants
- 755 R01MH116281-01A1 (BSA), R01ES030993-01A1 (LBA), and R01ES032547 (LBA). This
- 756 work was also supported by Singapore National Medical Research Council grants
- 757 NMRC/CIRG/1422/2015 and MOH-000032/MOHCIRG18may-0004 and the Singapore Ministry

758 of Health via the Duke-NUS Signature Research Programmes. LBA is an Abeloff V Scholar and

759 he is supported by an Alfred P. Sloan Research Fellowship. Research at UC San Diego was also

supported by a Packard Fellowship for Science and Engineering to LBA. AJG was funded by a

761 postdoctoral fellowship (grant nr. P2BSP3_178591). NP receives funding through the Cancer

762 Research UK Clinician Scientist Fellowship scheme and is supported by University College

763 London Cancer Institute. Research at Los Alamos National Laboratory was conducted under

764 Contract No. 89233218CNA000001 by the U.S. Department of Energy's National Nuclear

765 Security Administration and supported by Laboratory Directed Research and Development

766 (LDRD) grant 20190020DR (BSA). CDS is supported by the GEM consortium and

acknowledges funding for this work through a Cancer Research UK travel grant. The funders

had no roles in study design, data collection and analysis, decision to publish, or preparation of

the manuscript.

770

771

7	7	2
1	1	L

773 AUTHOR CONTRIBUTIONS

- T74 LBA and SMAI designed both SigProfilerExtractor's methodology and the performed analyses
- with help from NP, JZ, DJA, IM, BSA, LH, DCW, MTL, PB, MRS, and SGR. SMAI developed
- 776 SigProfilerExtractor with help from MV, ENB, YH, CDS, RV, and JW. All synthetic
- benchmarking datasets were generated by YW and SGR. SMAI documented
- 778 SigProfilerExtractor and performed the benchmarking of all tools on synthetic data with help
- from MDG, MB, BO, AK, and AA. Additional validations, confirmations, and applications of
- 780 SigProfilerExtractor to real and synthetic datasets were performed by SMAI, SM, SS, YRL, NS,
- 781 LR, TZ, AJG, YH, CDS, and SWB. LBA directed the overall research and wrote the manuscript
- 782 with help from SMAI and input from all other authors. All authors read and approved the final
- 783 manuscript.
- 784

785 COMPETING INTERESTS

MV is an employee of NVIDIA corporation. BSA and LBA are inventors of a US Patent
10,776,718 for source identification by non-negative matrix factorization. All other authors declare
no competing interests.

789

790 TOOL AVAILABILITY

791 SigProfilerExtractor and all its modules are open source and freely available for use under the

- permissive 2-clause BSD license. SigProfilerExtractor and its modules are implemented in
- 793 Python with an R wrapper package allowing users to run the tool from an R environment.
- 794 SigProfilerExtractor can be installed using the PyPI package manager from

- 795 <u>https://pypi.org/project/SigProfilerExtractor/</u> or downloaded from GitHub from
- 796 <u>https://github.com/AlexandrovLab/SigProfilerExtractor</u>. The R version of the tool can be
- 797 downloaded from <u>https://github.com/AlexandrovLab/SigProfilerExtractorR</u>. A detailed wiki
- 798 page including installation, usage, and explanation of result is provided at
- 799 <u>https://osf.io/t6j7u/wiki/home/</u>. SigProfilerExtractor is compatible with Windows, Linux, Unix,
- 800 and macOS operating systems.

801

802	REFF	CRENCE
803	1	Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R.
804		Deciphering signatures of mutational processes operative in human cancer. Cell Rep 3,
805		246-259, doi:10.1016/j.celrep.2012.12.008 (2013).
806	2	Bergstrom, E. N. et al. SigProfilerMatrixGenerator: a tool for visualizing and exploring
807		patterns of small mutational events. BMC Genomics 20, 685, doi:10.1186/s12864-019-
808		6041-2 (2019).
809	3	Lee, D. D. & Seung, H. S. Learning the parts of objects by non-negative matrix
810		factorization. Nature 401, 788-791, doi:10.1038/44565 (1999).
811	4	Févotte, C. & Cemgil, A. T. in 2009 17th European Signal Processing Conference.
812		1913-1917.
813	5	Dempster, A. P., Laird, N. M. & Rubin, D. B. Maximum Likelihood from Incomplete
814		Data Via the EM Algorithm. Journal of the Royal Statistical Society: Series B
815		(Methodological) 39 , 1-22, doi:10.1111/j.2517-6161.1977.tb01600.x (1977).
816	6	Suri, P. & Roy, N. R. in 2017 3rd International Conference on Computational
817		Intelligence & Communication Technology (CICT). 1-5.
818	7	Alexandrov, L. B. Understanding the origins of human cancer. Science 350, 1175,
819		doi:10.1126/science.aad7363 (2015).
820	8	Alexandrov, L. B. et al. Signatures of mutational processes in human cancer. Nature 500,
821		415-421, doi:10.1038/nature12477 (2013).
822	9	Petljak, M. & Alexandrov, L. B. Understanding mutagenesis through delineation of
823		mutational signatures in human cancer. Carcinogenesis 37, 531-540,
824		doi:10.1093/carcin/bgw055 (2016).
825	10	Pich, O. et al. The mutational footprints of cancer therapies. Nat Genet 51, 1732-1740,
826		doi:10.1038/s41588-019-0525-5 (2019).
827	11	Alexandrov, L. B. <i>et al.</i> Clock-like mutational processes in human somatic cells. <i>Nat</i>
828	10	<i>Genet</i> 47 , 1402-1407, doi:10.1038/ng.3441 (2015).
829	12	Alexandrov, L. B., Nik-Zainal, S., Siu, H. C., Leung, S. Y. & Stratton, M. R. A
830		mutational signature in gastric cancer suggests therapeutic strategies. <i>Nat Commun</i> 6 ,
831	12	8683, doi:10.1038/ncomms9683 (2015).
832	13	Alexandrov, L. B. <i>et al.</i> Mutational signatures associated with tobacco smoking in human
833	14	cancer. <i>Science</i> 354 , 618-622, doi:10.1126/science.aag0299 (2016).
834	14	Alexandrov, L. B. <i>et al.</i> The repertoire of mutational signatures in human cancer. <i>Nature</i>
835	15	578 , 94-101, doi:10.1038/s41586-020-1943-3 (2020).
836	13	Fischer, A., Illingworth, C. J., Campbell, P. J. & Mustonen, V. EMu: probabilistic
837		inference of mutational processes and their localization in the cancer genome. <i>Genome</i>
838	16	Biol 14, R39, doi:10.1186/gb-2013-14-4-r39 (2013).
839 840	10	Mayakonda, A., Lin, D. C., Assenov, Y., Plass, C. & Koeffler, H. P. Maftools: efficient and comprehensive analysis of somatic variants in cancer. <i>Genome Res</i> 28 , 1747-1756,
840 841		doi:10.1101/gr.239244.118 (2018).
842	17	Blokzijl, F., Janssen, R., van Boxtel, R. & Cuppen, E. MutationalPatterns: comprehensive
843	1 /	genome-wide analysis of mutational processes. <i>Genome Med</i> 10 , 33,
844		doi:10.1186/s13073-018-0539-0 (2018).
845	18	Fantini, D., Vidimar, V., Yu, Y., Condello, S. & Meeks, J. MutSignatures: an R package
846	10	for extraction and analysis of cancer mutational signatures. <i>Scientific Reports</i> 10 , 18217-
847		18217 (2020).
517		

848 19 Ardin, M. et al. MutSpec: a Galaxy toolbox for streamlined analyses of somatic mutation 849 spectra in human and mouse cancer genomes. BMC Bioinformatics 17, 170, 850 doi:10.1186/s12859-016-1011-z (2016). 851 20 Gori, K. & Baez-Ortega, A. sigfit: flexible Bayesian inference of mutational signatures. 852 bioRxiv, 372896, doi:10.1101/372896 (2020). 853 21 Wang, S. et al. Copy number signature analyses in prostate cancer reveal distinct 854 etiologies and clinical outcomes. medRxiv, 2020.2004.2027.20082404, 855 doi:10.1101/2020.04.27.20082404 (2020). 856 Kasar, S. et al. Whole-genome sequencing reveals activation-induced cytidine deaminase 22 857 signatures during indolent chronic lymphocytic leukaemia evolution. Nat Commun 6, 858 8866, doi:10.1038/ncomms9866 (2015). 859 Taylor-Weiner, A. et al. Scaling computational genomics to millions of individuals with 23 860 GPUs. Genome Biol 20, 228, doi:10.1186/s13059-019-1836-7 (2019). 861 24 Degasperi, A. et al. A practical framework and online tool for mutational signature 862 analyses show inter-tissue variation and driver dependencies. Nat Cancer 1, 249-263, 863 doi:10.1038/s43018-020-0027-5 (2020). 864 25 Rosales, R. A., Drummond, R. D., Valieris, R., Dias-Neto, E. & da Silva, I. T. signeR: an 865 empirical Bayesian approach to mutational signature discovery. Bioinformatics 33, 8-16, doi:10.1093/bioinformatics/btw572 (2017). 866 Gehring, J. S., Fischer, B., Lawrence, M. & Huber, W. SomaticSignatures: inferring 867 26 868 mutational signatures from single-nucleotide variants. *Bioinformatics* **31**, 3673-3675, 869 doi:10.1093/bioinformatics/btv408 (2015). 870 27 Vöhringer, H. & Gerstung, M. Learning mutational signatures and their multidimensional 871 genomic properties with TensorSignatures. bioRxiv, 850453 (2019). 872 28 Consortium, I. T. P.-C. A. o. W. G. Pan-cancer analysis of whole genomes. Nature 578, 873 82-93, doi:10.1038/s41586-020-1969-6 (2020). 874 29 Kullback, S. & Leibler, R. A. On Information and Sufficiency. Ann. Math. Statist. 22, 79-86, doi:10.1214/aoms/1177729694 (1951). 875 876 30 Kuhn, H. W. The Hungarian method for the assignment problem. Naval research 877 logistics quarterly 2, 83-97 (1955). 878 Huang, K., Sidiropoulos, N. D. & Swami, A. Non-Negative Matrix Factorization 31 879 Revisited: Uniqueness and Algorithm for Symmetric Decomposition. IEEE Transactions 880 on Signal Processing 62, 211-224, doi:10.1109/TSP.2013.2285514 (2014). 881 32 Lin, C. On the Convergence of Multiplicative Update Algorithms for Nonnegative Matrix 882 Factorization. IEEE Transactions on Neural Networks 18, 1589-1596, 883 doi:10.1109/TNN.2007.895831 (2007). 884 Bergstrom, E. N., Barnes, M., Martincorena, I. & Alexandrov, L. B. Generating realistic 33 885 null hypothesis of cancer mutational landscapes using SigProfilerSimulator. BMC 886 Bioinformatics 21, 438, doi:10.1186/s12859-020-03772-3 (2020). 887 34 Tan, V. Y. F. & Févotte, C. Automatic Relevance Determination in Nonnegative Matrix 888 Factorization with the /spl beta/-Divergence. IEEE Transactions on Pattern Analysis and 889 Machine Intelligence 35, 1592-1605, doi:10.1109/TPAMI.2012.240 (2013). 890 35 Lawson, A. R. J. et al. Extensive heterogeneity in somatic mutation and selection in the 891 human bladder. Science 370, 75-82, doi:10.1126/science.aba8347 (2020).

892 36 Moody, S. et al. Mutational signatures in esophageal squamous cell carcinoma from eight 893 countries of varying incidence. medRxiv, 2021.2004.2029.21255920, 894 doi:10.1101/2021.04.29.21255920 (2021). 895 37 Gaujoux, R. & Seoighe, C. A flexible R package for nonnegative matrix factorization. 896 BMC Bioinformatics 11, 367, doi:10.1186/1471-2105-11-367 (2010). 897 Brunet, J. P., Tamayo, P., Golub, T. R. & Mesirov, J. P. Metagenes and molecular pattern 38 898 discovery using matrix factorization. Proc Natl Acad Sci USA 101, 4164-4169, 899 doi:10.1073/pnas.0308531101 (2004). 900 Abadi, M. et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous 39 901 Distributed Systems. arXiv e-prints, arXiv:1603.04467 (2016). 902 Carpenter, B. et al. Stan: A Probabilistic Programming Language. Journal of Statistical 40 903 Software 76, doi:10.18637/jss.v076.i01 (2017). 904 Schwarz, G. Estimating the Dimension of a Model. The Annals of Statistics 6, 461-464, 41 905 doi:10.1214/aos/1176344136 (1978). Thorndike, R. L. Who belongs in the family? Psychometrika 18, 267-276, 906 42 907 doi:10.1007/bf02289263 (1953). 908 43 Benjamin, N., Raviteja, V., Miguel, A. H.-H., Svetlana, K. & Boian, A. A neural network 909 for determination of latent dimensionality in Nonnegative Matrix Factorization. Machine 910 Learning: Science and Technology (2020). 911 44 Kim, J. et al. Somatic ERCC2 mutations are associated with a distinct genomic signature 912 in urothelial tumors. Nat Genet 48, 600-606, doi:10.1038/ng.3557 (2016). 913 Cancer Genome Atlas Research, N. Comprehensive molecular characterization of 45 914 urothelial bladder carcinoma. Nature 507, 315-322, doi:10.1038/nature12965 (2014). 915 46 Shalabi, L. A. & Shaaban, Z. in 2006 International Conference on Dependability of 916 Computer Systems. 207-214. 917 47 Žitnik, M. & Zupan, B. Nimfa: A python library for nonnegative matrix factorization. 918 The Journal of Machine Learning Research 13, 849-853 (2012). 919 Lew, J. et al. in 2019 IEEE International Symposium on Performance Analysis of Systems 48 920 and Software (ISPASS). 151-152. 921 Aranganayagi, S. & Thangavel, K. in International Conference on Computational 49 922 Intelligence and Multimedia Applications (ICCIMA 2007). 13-17 (IEEE). 923 50 Franc, V., Hlaváč, V. & Navara, M. in Computer Analysis of Images and Patterns. (eds 924 André Gagalowicz & Wilfried Philips) 407-414 (Springer Berlin Heidelberg). 925 51 Stacklies, W., Redestig, H., Scholz, M., Walther, D. & Selbig, J. pcaMethods--a 926 bioconductor package providing PCA methods for incomplete data. Bioinformatics 23, 927 1164-1167, doi:10.1093/bioinformatics/btm069 (2007). 928

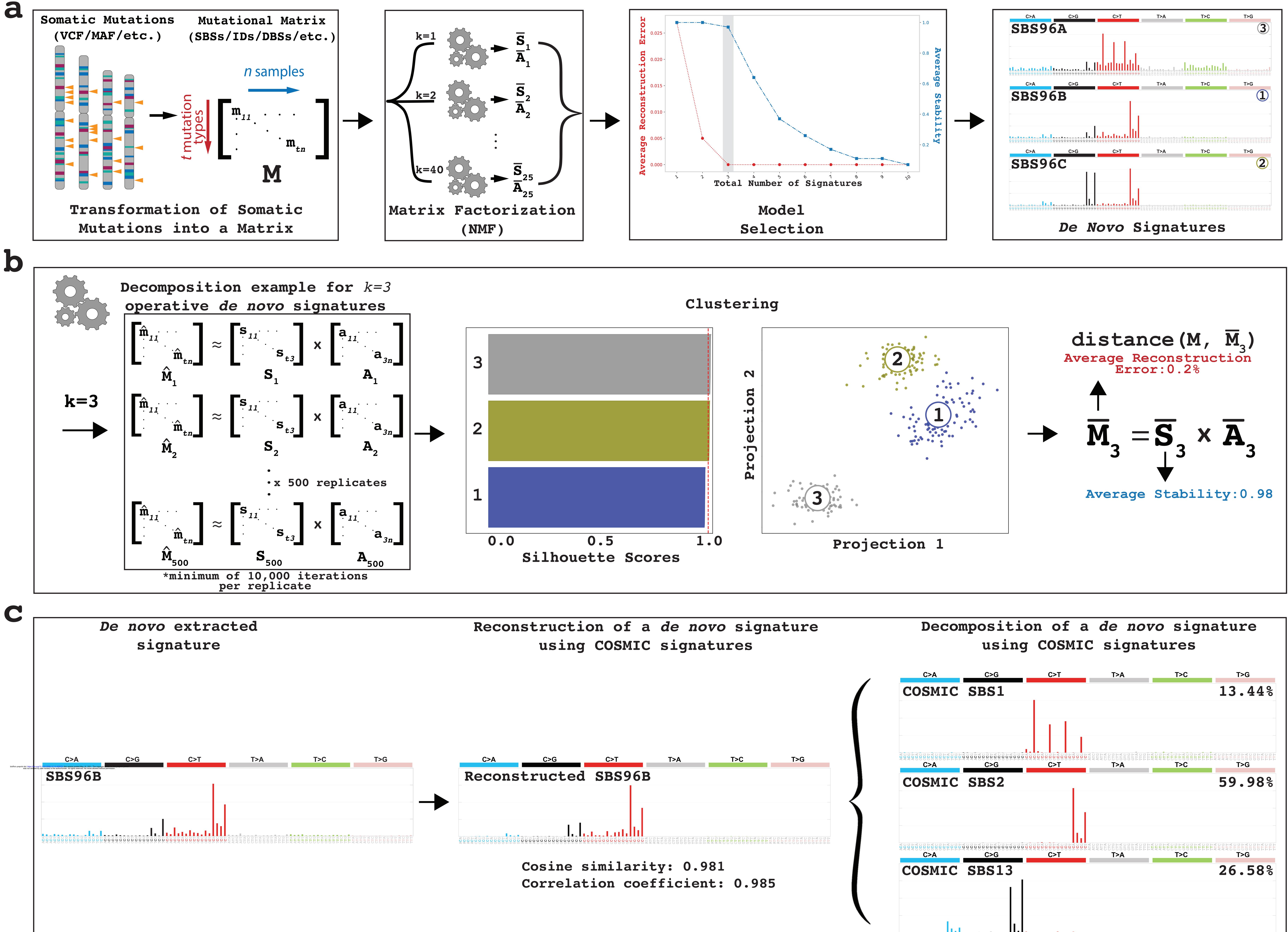
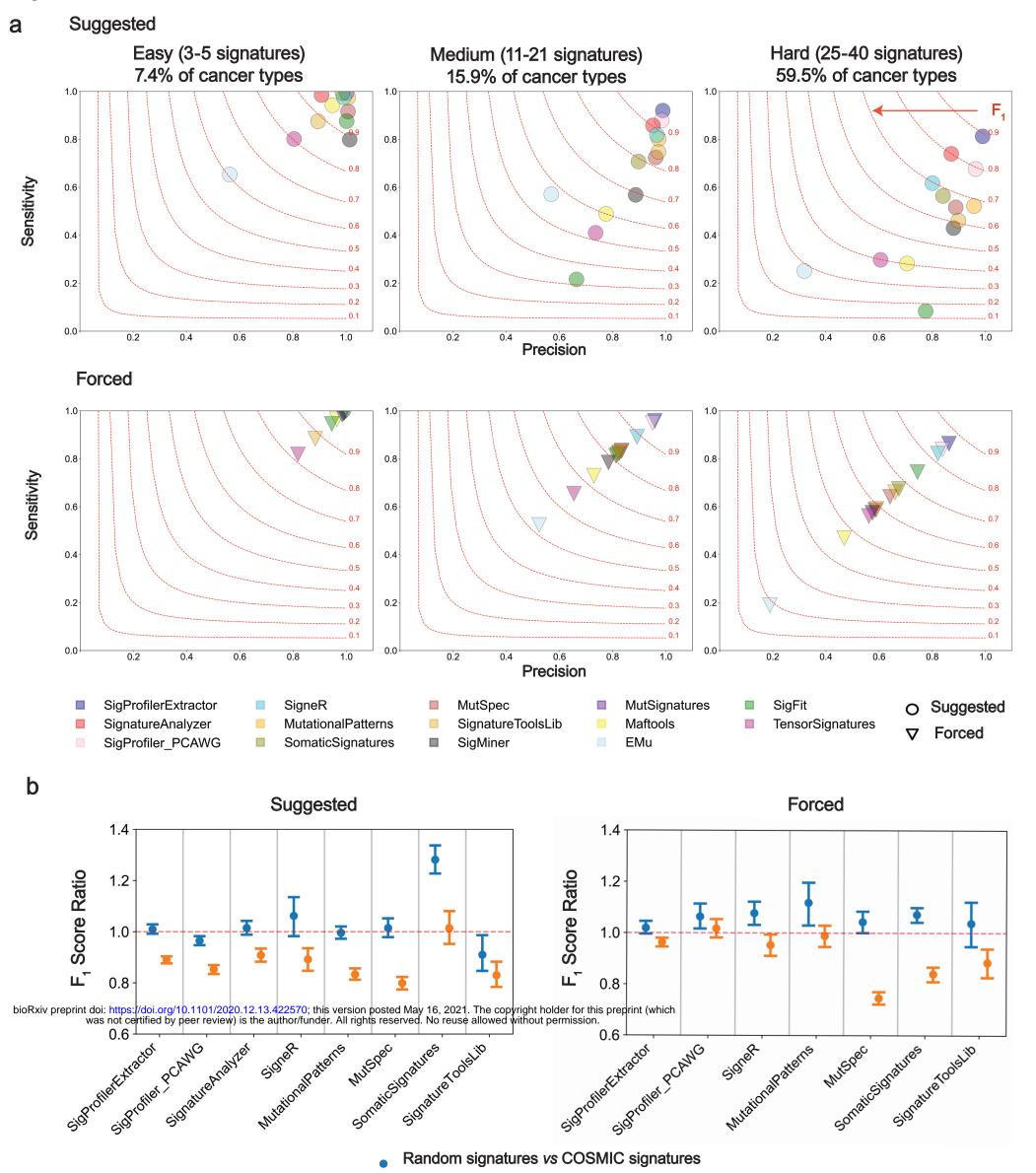
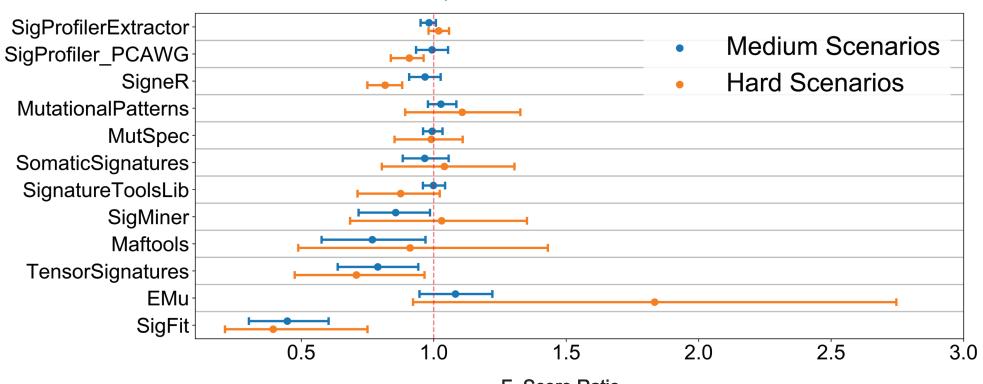


Fig.2

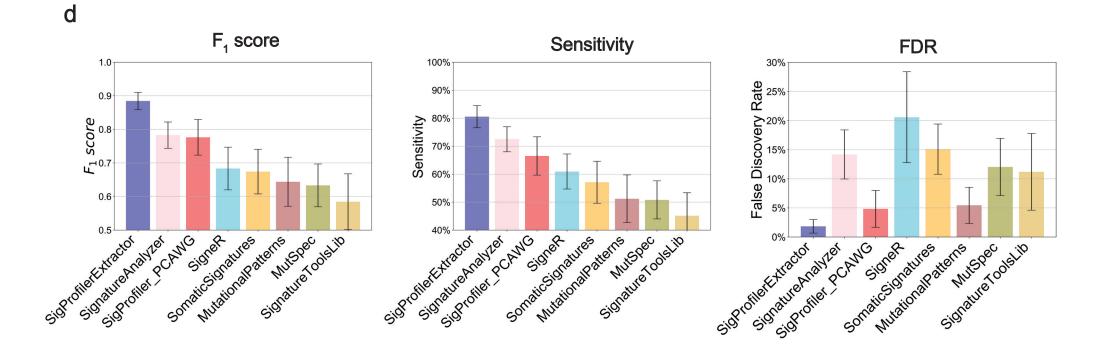


- Random signatures vs COSMIC signatures
- SA signatures vs COSMIC signatures

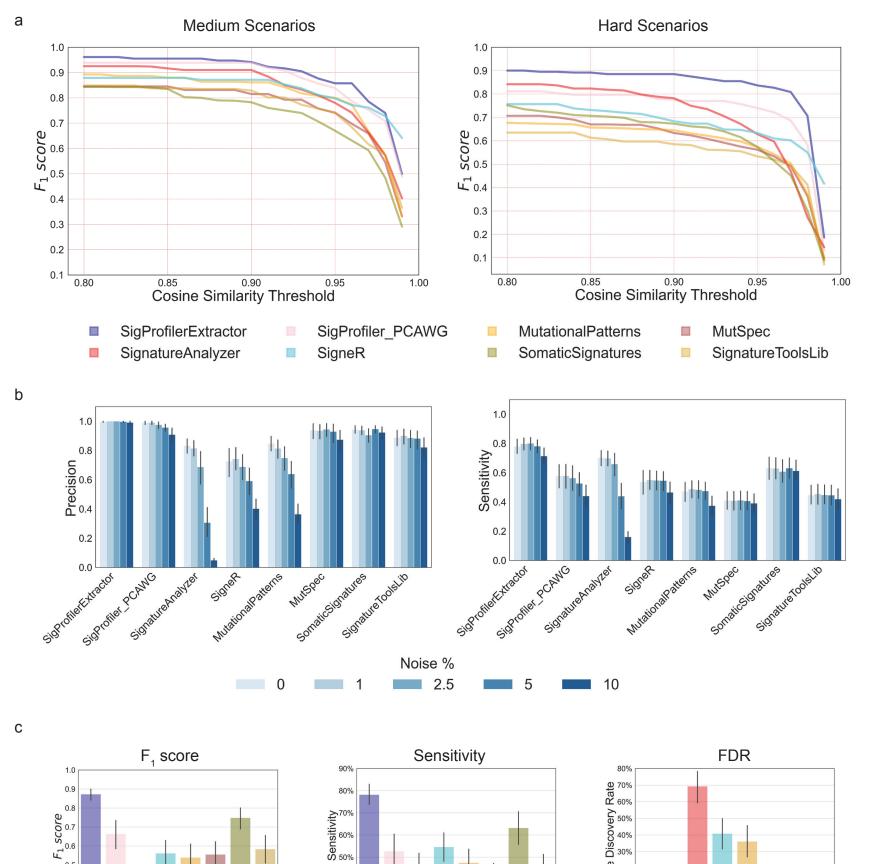


F₁ Score Ratio

F₁ Score (Suggested/Forced Solutions)



С



40%

Frunder PCANG

Signature Analyter

Signofie Extractor

Wulaional atens

SOR THE TOOLED

Sonalcoolatures.

0.5

0.4

Frunder PCANG

Signature Analyter

signoner charter

MulalionalPatients

Signature Tobship

Sonalcoolatures.

Ealse 10%

Signalian and a company of the second

Funder Dear Provent

und inthe halfer

Nutational Patients

Sonalosonalues Signatue Tobbility

