

HiLDA: a statistical approach to investigate differences in mutational signatures

Zhi Yang¹, Priyatama Pandey¹, Darryl Shibata², David V. Conti¹, Paul Marjoram¹, and Kimberly D. Siegmund¹

¹Department of Preventive Medicine, Keck School of Medicine of the University of Southern California, Los Angeles, CA USA

²Department of Pathology, Keck School of Medicine of the University of Southern California, Los Angeles, CA USA

Corresponding author:

Zhi Yang¹

Email address: zhiyang@usc.edu

ABSTRACT

We propose a hierarchical latent Dirichlet allocation model (HiLDA) for characterizing somatic mutation data in cancer. The method allows us to infer mutational patterns and their relative frequencies in a set of tumor mutational catalogs and to compare the estimated frequencies between tumor sets. We apply our method to somatic mutations in colon cancer with mutations classified by the time of occurrence, before or after tumor initiation. Applying the methods to 16 colon cancers, we found significant associations between the relative frequencies of mutational patterns and the time of occurrence of mutations. Our novel method provides higher statistical power for detecting differences in mutational signatures.

INTRODUCTION

A variety of mutational processes occur over the lifetime of an individual, and thereby uniquely contribute to the catalog of somatic mutations observed in a tumor. Some processes leave a molecular signature: a specific base substitution occurring within a particular pattern of neighboring bases. A variety of methods exist to discover mutational signatures from the catalog of all somatic mutations in a set of tumors, estimating the latent mutational signatures as well as the latent *exposures* (i.e., fraction of mutations) each signature contributes to the total catalog. The first large study of mutational signatures in cancer identified variation in mutational signatures and mutational exposures across 21 different cancer types (Alexandrov et al., 2013). To better understand the sources of variation in the mutational exposures across cancers, our interest is in statistical methods used to characterize these latent mutational exposures across different cancer subtypes. Moreover, by classifying mutations by their time of occurrence, before or after tumor initiation, we can investigate whether new mutational processes occur during tumor growth.

Previous studies interested in comparing mutational exposure estimates between different groups of tumor catalogs conducted a post hoc analysis. The analysis proceeded in two stages. First, they performed one of the several different approaches for mathematically extracting the latent mutational signatures and their exposures from the mutational catalogs (see Baez-Ortega and Gori (2017) for a review of such methods). Later, they conducted an independent test of association between the point estimates of the mutational exposures and external covariates. Examples of covariates included cancer subtype, or patient history of alcohol or tobacco use. A common choice for the second stage test is a Wilcoxon rank-sum test (Mann and Whitney, 1947; Network et al., 2017; Chang et al., 2017; Hillman et al., 2017; Letouzé et al., 2017; Meier et al., 2018; Haradhvala et al., 2018; Qin et al., 2018; Olivier et al., 2019; Guo et al., 2018). However, the variation of the exposure estimates is affected by two factors, the number of mutations in the tumor and the variation in exposure frequency in the patient population. The former, the number of mutations in the tumor, affects the accuracy of the exposure estimates. The application of the Wilcoxon rank-sum test on the exposure estimates does not take into consideration their accuracy, which can lead to loss of efficiency and test power. We address this by introducing a unified parametric model for testing variation of mutational exposures between groups of mutational catalogs, where the exposure frequencies

47 are modeled using a Dirichlet distribution.

48 We propose a hierarchical latent Dirichlet allocation model (HiLDA) that adds an additional level to the
49 latent Dirichlet allocation (LDA) model from Shiraishi et al. (2015). Shiraishi's model, like the majority
50 of deconvolution approaches, focuses on signatures for single-nucleotide substitutions, characterizing the
51 mutation types by context, using local features in the genome such as the pattern of flanking bases and
52 possibly the transcription strand. For both model parsimony and interpretation, we choose to extend their
53 LDA model. First, it requires fewer parameters than competing methods, giving it higher power to detect
54 patterns 5 bases in length compared to other models that consider only 3-base contexts (Shiraishi et al.,
55 2015). Second, signature visualization methods lead to easy interpretation; an example is the common
56 C>T substitution at CpG sites instead of the more complicated NpCpG patterns that appear when using
57 the trinucleotide context. Like the LDA model, HiLDA retains all the functionality for estimating both
58 the latent signatures and the latent mutational exposure of each signature for each tumor catalog. Our
59 newly-added hierarchical level allows HiLDA to simultaneously test whether those mean exposures differ
60 between different groups of catalogs while accounting for the uncertainty in the exposure estimates.
61 Additionally, we can now parse out differences in group means in the presence of differences in group
62 variances, which is not tenable when using post hoc nonparametric location-scale tests.

63 In this paper, we use HiLDA to study the association between the mutational exposures and the time
64 of mutation occurrence in tumorigenesis. We classify cancer mutations into *trunk* or *branch* mutations:
65 trunk mutations being those that occur before growth of the tumor, while branch mutations are those
66 that occur during the tumor expansion process. A test of whether mutational exposures differ by time
67 of mutation occurrence will allow us to assess whether new mutational processes occur following the
68 transformation of the first cancer cell.

69 METHODS

70 Hierarchical Bayesian Mixture Model

71 We introduce a hierarchical latent Dirichlet allocation model (HiLDA) using the following notation,
72 also summarized in Table 1. Let i index the mutational catalog and j the mutation. The nucleotide
73 substitutions are reduced to six possible types (C>A, C>T, C>G, T>A, T>C, T>G) to eliminate
74 redundancy introduced by the complementary strands. Each observed mutation is characterized by
75 a vector, $\mathbf{X}_{i,j}$ describing the nucleotide substitution (e.g. C>T) and a set of genomic features in the
76 neighborhood. Example features include the base(s) 3' and 5' of the nucleotide substitution (C, G, A, T),
77 and the transcription strand (+, -). Each observed feature characteristic, $x_{i,j,l}$ for mutation feature l , takes
78 values in the set $\{1, 2, \dots, M_l\}$ (where $M_l = 6$ for the nucleotide substitution, or 4 for a flanking base, and
79 2 for the transcription strand).

80 We assume each mutation belongs to one of K distinct signatures. A specific mutational signature k is
81 defined by an l -tuple of probability vectors, \mathbf{F}_k , denoting the relative frequencies of the M_l discrete values
82 for the l features, i.e., a vector $\mathbf{f}_{k,l}$ for the M_l values corresponding to feature l . We let $z_{i,j}$ denote the unique
83 latent assignment of mutation $\mathbf{X}_{i,j}$ to a particular signature. Then, given the signature to which a mutation
84 belongs, the probability of observing a mutational pattern is calculated as the product of the mutation
85 feature probabilities for that signature. Thus, for signature k we write $Pr(\mathbf{X}_{i,j}|z_{i,j}) = \prod_l f_{k,l}(x_{i,j,l}|z_{i,j})$.
86 This assumes independent contributions of each feature to the signature. To model each multinomial
87 distribution of $\mathbf{f}_{k,l}$, we use a non-informative Dirichlet prior distribution with all concentration parameters
88 equal to one.

89 The unique personal exposure history of each individual leads to them having a particular (latent)
90 vector, \mathbf{q}_i , indicating the resulting contribution of each of the K signatures to that individual's mutational
91 catalog. These \mathbf{q}_i s are modeled using a Dirichlet distribution with concentration parameters $\boldsymbol{\alpha}$, i.e.,
92 $\mathbf{q}_i \sim Dir(\boldsymbol{\alpha})$. Extending this model to the two-group setting, we allow the Dirichlet parameters to depend
93 on group, $Dir(\boldsymbol{\alpha}^{(g_i)})$, with g_i indexing the group corresponding to the i th catalog ($g_i = 1$ or 2). The mean
94 mutational exposures, $E(\mathbf{q}_i)$, denoted by $\boldsymbol{\mu}^{(g_i)}$, are represented by using the concentration parameters, i.e.,
95 $\boldsymbol{\mu}^{(g_i)} = \boldsymbol{\alpha}^{(g_i)} / \sum \boldsymbol{\alpha}^{(g_i)}$.

96 With this extension, we can infer differences in mutational processes between groups of catalogs by
97 testing whether the mean mutational exposures differ between the two sets, i.e., at least one $\mu_k^{(1)} \neq \mu_k^{(2)}$.
98 The likelihood and prior of the multi-level model is specified as follows,

$$\begin{aligned}x_{i,j,l}|z_{i,j} &\sim \text{Multinomial}(\mathbf{f}_{z_{i,j},l}) \\z_{i,j} &\sim \text{Multinomial}(\mathbf{q}_i|g) \\ \mathbf{q}_i|g_i &\sim \text{Dir}(\boldsymbol{\alpha}^{(g_i)})\end{aligned}$$

99 For full details see See Text S1. and Fig. S2..

100 Testing for Differences in Signature Exposures

101 To characterize the signature contributions for different sets of tumor catalogs, we wish to conduct a
102 hypothesis test that there is no difference in mean exposures versus the alternative that the mean exposure
103 of at least one signature differs between the two groups, i.e. $H_0 : \boldsymbol{\mu}^{(1)} = \boldsymbol{\mu}^{(2)}$ vs. H_1 : at least one
104 $\mu_k^{(1)} \neq \mu_k^{(2)}$. We propose both local and global tests, implemented in a Bayesian framework. The former
105 provides signature-level evaluations to determine where the differences in mean mutational exposures
106 occur, while the latter provides an overall conclusion about any difference in mean mutational exposures.
107 The details of our implementation are given in our Just Another Gibbs Sampler (JAGS) scripts and Source
108 code is freely available in Github at <https://github.com/USCbiostats/HiLDA> (Plummer et al., 2003).

109 A local test to identify signatures with different exposures

We propose a signature-level (local) hypothesis test to allow us to infer which signature(s) contribute a
different mean exposure to the mutational catalogs across tumor sets, i.e., $\mu_k^{(1)} \neq \mu_k^{(2)}$. To measure the
difference between mean signature exposure vectors, we implement HiLDA by specifying two Dirichlet
distributions, $\text{Dir}(\boldsymbol{\alpha}^{(1)})$ and $\text{Dir}(\boldsymbol{\alpha}^{(2)})$, as priors for the distribution of mutational exposures \mathbf{q}_i of each
group (Spiegelhalter et al., 2003). Using this formulation, the difference between the two groups of the
mean exposure of signature k is calculated as,

$$\Delta_k = \mu_k^{(2)} - \mu_k^{(1)} = \frac{\alpha_k^{(2)}}{\sum_k \alpha_k^{(2)}} - \frac{\alpha_k^{(1)}}{\sum_k \alpha_k^{(1)}} \quad (1)$$

For all parameters, $\alpha_k^{(1)}$'s and $\alpha_k^{(2)}$'s, we use independent, non-informative gamma distribution priors
with a rate of 0.001 and shape of 0.001; this results in a mean of 1 and variance of 1000. So,

$$\alpha_k^{(g_i)} \sim \text{Gamma}(0.001, 0.001)$$

110 We estimate parameters via Markov chain Monte Carlo (MCMC) using two chains (Carlin and Chib,
111 1995). We assess convergence of the two MCMC chains using the potential scale reduction factor (Rhat)
112 in Gelman et al. (1992), which is required to be less than or equal to 1.05 for all parameters in order to
113 conclude that the MCMC run has converged. After obtaining the posterior distribution of the differences
114 (i.e., of Δ_k), there are two possible approaches to performing inference. We can: 1) use the Wald test to
115 compute the P -value using the means and standard errors of the posterior distribution for Δ_k ; 2) determine
116 whether the 95% credible interval of the posterior distribution for Δ_k contains zero.

117 A global test using the Bayes factor

We also propose a global test to provide an overall conclusion on whether the mean exposures differ
between groups of catalogs. It uses the Bayes factor, the ratio of posterior to prior odds in favor of the
alternative (H_1 : at least one $\mu_k^{(1)} \neq \mu_k^{(2)}$, $k = 1, \dots, K$) compared to the null (H_0 : $\boldsymbol{\mu}^{(1)} = \boldsymbol{\mu}^{(2)}$), to indicate
the strength of evidence that they do differ, without explicit details on how they differ. Thus, we can
calculate the Bayes factor as:

$$\text{Bayes Factor} = \frac{\text{Pr}(H_1|Data)}{\text{Pr}(H_0|Data)} \bigg/ \frac{\text{Pr}(H_1)}{\text{Pr}(H_0)}. \quad (2)$$

118 Since the likelihood is analytically intractable, the Bayes factor is calculated via MCMC (Carlin and
119 Chib, 1995). In order to estimate the Bayes factor, during the MCMC analysis, a single binary hypothesis
120 index variable is used to indicate which hypothesis explains the observed data (Lodewyckx et al., 2011).
121 The parameters of two Dirichlet distributions, $\text{Dir}(\boldsymbol{\alpha}^{(1)})$ and $\text{Dir}(\boldsymbol{\alpha}^{(2)})$, are drawn from the same prior if

122 the index takes the value 1, whereas they are drawn from different priors if it takes the value 2. Initially,
123 the prior hypothesis odds is set to be $0.5/0.5 = 1$, which means that both hypotheses are assumed equally
124 likely under the prior. In order to improve computational efficiency in extreme situations in which one
125 hypothesis dominates the other, we can use a different prior odds value (Carlin and Chib, 1995).

126 **Two-stage Inference Methods using the Point Estimates of mutational exposures**

127 An alternative approach is to perform hypothesis testing using point estimates of the mutational exposures,
128 \hat{q}_i , in a two-stage analysis, which we refer to as the "two-stage" method (TS). We used the R package
129 **pmsignature** to estimate \hat{q} (Shiraishi et al., 2015). Other methods are also available, but we selected
130 **pmsignature** for the purpose of comparisons to the results from HiLDA since it assumes the same model
131 for estimating signatures under independence of features. We summarize the steps of the TS method as
132 follows:

- 133 1. Jointly estimate the vectors of mutational signature exposures, q_i , for each mutational catalog.
- 134 2. Test for differential mutational exposures for signature k by performing the Wilcoxon rank-sum test
135 on the \hat{q}_k .

136 However, we note that the Wilcoxon rank-sum test in stage 2 is also sensitive to changes in variance
137 across the two groups, which might lead to significant results even when there has been no change in
138 mean exposures (Kasuya, 2001; Ruxton, 2006). We implemented the two-stage method using R version
139 3.5.0 (R Core Team, 2017). A two-sided P value of less than 0.05 was considered statistically significant.

140 **Choosing the Number of Signatures**

141 The number of signatures, K , needs to be determined prior to any of the above analyses. We adopted the
142 method of Shiraishi et al. (2015) to determine K . Their method is based on the following criteria:

- 143 1. The optimal value of K is selected over a range of K values such that the likelihood remains
144 relatively high while simultaneously having relatively low standard errors for the parameters.
- 145 2. Pairwise correlations between any two signatures (the k th signature and the k' th signature, say)
146 are measured by calculating the Pearson correlation between their estimated exposures across all
147 samples, (i.e., the correlation between $(\hat{q}_{1,k}, \dots, \hat{q}_{L,k})$ and $(\hat{q}_{1,k'}, \dots, \hat{q}_{L,k'})$). K is chosen such that no
148 strong correlation (i.e., >0.6) exists between any pair.

149 For full details see Shiraishi et al. (2015).

150 **Application to Tumor Evolution**

151 ***USC Colon Cancer Data***

152 Our goal is to identify whether any new mutational signatures occur during colon cancer growth that
153 distinguish cancer evolution from normal tissue evolution. To achieve this, we classify somatic mutations
154 into two catalogs according to time of occurrence: those that accumulated between the time of the zygote
155 and the first tumor cell, which we call trunk mutations, and those that occur *de novo* during tumor growth,
156 which we refer to as branch mutations. We then estimate mutational signatures in the two sets of catalogs
157 and test whether the mean mutational exposures differ between them.

158 We analyzed a total of 16 colon tumors. Tumor and adjacent normal tissue were subject to whole
159 exome sequencing, and somatic mutations called using the GATK pipeline and MuTect (details below).
160 Somatic mutations in the tumors were defined as nucleotide variants that were detected in tumor tissue
161 but did not also appear in the patient-matched normal tissue. We used multi-region tumor sampling to
162 allow us to distinguish between trunk from branch mutations (Siegmund and Shibata, 2016). Each tumor
163 was sampled twice, with bulk tissue samples taken from opposite tumor halves. We classified somatic
164 mutations appearing in both tumor halves as trunk, because only trunk mutations are likely to appear
165 in both tumor halves, while mutations found on only one side of a tumor were labeled as branch. This
166 approach has previously been shown to be 99% sensitive for calling trunk mutations and 85% sensitive
167 for calling branch mutations (Siegmund and Shibata, 2016). Fifteen of the 16 tumors were previously
168 analyzed in a study of cell motility (Ryser et al., 2018).

169 The sequence data were processed using the GATK pipeline version 3.7 (DePristo et al., 2011) and
170 somatic mutations called with MuTect version 1.1.7 (Cibulskis et al., 2013), applying the quality filters
171 KEEP (default parameters) and COVERED (read depth of 14 in tumor and 10 in matched normal - use of

172 a lower coverage threshold in normal tissue is as recommended in (Cibulskis et al., 2013)). We excluded
173 any mutations that either had an allele frequency less than 0.10, because sequencing errors are more
174 common among low-frequency mutations (Cibulskis et al., 2013), or that were not also found by Strelka
175 (Saunders et al., 2012), which we used as a confirmatory control. Somatic mutations on chromosomes 1
176 to 22 were used for mutational signature analysis.

177 RESULTS

178 Application to Tumor Evolution

179 A total of 12,554 somatic single-nucleotide substitutions were identified, with a median of 277 per sample
180 (range: 82 - 1,762) (See Table S3.). One tumor with microsatellite instability has more than double the
181 number of somatic mutations (1751 side A, 1762 side B) than any of the remaining 30 catalogs (all <750
182 mutations). In our first analysis, we compared the mutational exposures in side A to those in side B. If the
183 tumors represent a single clonal expansion, we would expect similar mutational exposure frequencies in
184 the two catalogs from the same tumor. Indeed, this is what we found (Table 2).

185 We identified a median of 174 trunk and 186 branch mutations per tumor. The numbers ranged from 49
186 to 1,578 trunk mutations and from 66 to 503 branch mutations (Fig. 1A). Interestingly, the microsatellite
187 unstable tumor had the most trunk mutations, but not the most branch mutations, suggesting that during
188 tumor growth the mutation frequency is similar in microsatellite stable and unstable tumors. Fig. 1B
189 shows that the C>T substitution is most common in all trunk catalogs, and most branch catalogs. The
190 spontaneous deamination of methylated Cs in CpGs is known to contribute to hotspots of C>T mutation
191 in the genome.

192 We identified three mutational signatures in our data (see Fig. S4.). Those three signatures, and their
193 corresponding exposures, are depicted in Fig. 2. The signature shown in the yellow box in the same
194 figure, involving C>T mutations at NpCpG sites, resembles signature 7 in Shiraishi et al. (2015), where
195 it was identified in 25 out of 30 cancer types and likely relates to the deamination of 5-methylcytosine
196 ('aging'); the signature in the orange box, involving T>G mutations at GpGpTpGpN sites, is novel; the
197 third signature, in the red box, is qualitatively similar to signature 17 in Shiraishi et al. (2015), reflecting a
198 signal specific to colorectal cancers. The pairwise cosine similarities between pairs of signatures are 0.12,
199 0.01, and 0.02 which are rather dissimilar from each other given the [0, 1] range for cosine similarity.
200 Using HiLDA, we test whether the three signatures differ in mean exposure between trunk and branch
201 mutations.

202 Our global test strongly suggests that, in our data, the signature exposures statistically differ between
203 trunk and branch catalogs (Bayes Factor 1265.0). Each of the individual signatures (depicted in Fig. 2B) is
204 found to differ in exposure between the two sample groups, a conclusion supported by both HiLDA and the
205 two-stage method (Table 3). From Fig. 2A, it is evident that the exposures of the first ('aging') signature
206 in trunk mutations is almost always greater than that for the matching catalog of branch mutations, which
207 is intuitively consistent with the fact that trunk mutations may well reflect an accumulation of mutations
208 over the life of the subject, whereas branch mutations are accumulated only after tumor initiation. For the
209 previously unseen signature, the higher exposures in branch catalogs might suggest that this signature's
210 underlying mechanism for generating mutations might be associated with the processes occurring during
211 tumor evolution as opposed to normal development. From Fig. 2C, we observed that the distributional
212 ranges of the two groups of mutational exposures have some overlaps, but that the centers of each group,
213 i.e., the means of mutational exposures, are clearly deviated from each other. However, the distributional
214 radii, indicating the variances of mutational exposures, do not substantially differ between the groups.

215 We sought to validate the discovery of the previously unseen signature using both targeted sequencing
216 data from the same tumor set (Siegmond and Shibata, 2016) and using publicly available data from the
217 Cancer Genome Atlas. Four T>G substitutions that we assigned to the previously unseen signature
218 were part of an independent validation set of mutations subjected to targeted, high-coverage Ampliseq
219 technology (Siegmond and Shibata, 2016); all four of these T>G substitutions failed to validate. Further,
220 a systematic analysis of data from the Cancer Genome Atlas Williams et al. (2016) also did not find
221 evidence for this signature. Therefore, we cannot rule out that the signature is the result of sequencing
222 error. We now go on to assess the reliability of results using a simulation study.

223 Simulation Study

224 We conducted a simulation study to assess the performance of both HiLDA and the two-stage approach in
225 terms of the false-positive rate (FPR) and true-positive rate (TPR), in local, univariate tests of the difference
226 in mean exposure between two groups of mutational catalogs. In order to assess the functionality of the
227 methods in a setting similar to that of the USC data, we simulate somatic mutations directly using the
228 estimated signatures (\mathbf{f}_k) from Fig. 2 for the same number of mutational catalogs (two groups of 16
229 catalogs each) and somatic mutations per catalog (J_i in S3 Table). The mutational exposures (\mathbf{q}_i) were
230 indirectly used to derive the concentration parameters of the Dirichlet distributions. The scenarios are as
231 follows:

- 232 1. The two groups of mutational catalogs are from separate Dirichlet distributions with param-
233 eters $\boldsymbol{\alpha}^{(1)} = (9.2, 0.2, 7.5)$ and $\boldsymbol{\alpha}^{(2)} = (4.2, 0.6, 7.3)$. Here, the $\boldsymbol{\alpha}$ s corresponds to the maximum-
234 likelihood estimated parameters from the three exposure distributions in the trunk and branch muta-
235 tional catalogs. This gives mean exposures of $\boldsymbol{\mu}^{(1)} = (0.54, 0.01, 0.44)$ and $\boldsymbol{\mu}^{(2)} = (0.35, 0.05, 0.60)$
236 in trunk and branch catalogs, respectively, for the aging signature, new signature, and random
237 signature.
- 238 2. The two groups of mutational catalogs are from the same Dirichlet distribution, $Dir(4.2, 0.6, 7.3)$,
239 (so here we use the concentration parameters estimated from the branch mutational catalogs).

240 For each tumor, mutational exposures \mathbf{q}_i , are drawn from the Dirichlet distribution. Each set of
241 probabilities parameterize a multinomial distribution later used to probabilistically choose the underlying
242 mutational signature for a mutation (See Fig. S5.). Then, every mutation feature in the mutational
243 pattern of the mutation is simulated independently from a corresponding multinomial distribution of the
244 chosen signature. To estimate the FPRs, 1000 sets of data were simulated for scenario 2, when there is
245 no difference in the exposure distribution between two groups of mutational catalogs. The two-stage
246 method is slightly conservative for 1st and 3rd signatures (resulting FPRs of 4.3%, 5.2%, and 4.3%) when
247 testing at the 5% significant level (Table 4). In comparison, HiLDA showed better control of the FPR by
248 using the 95% credible interval of the posterior distributions (4.8%, 5.0%, and 5.1%). The Wald test also
249 showed control of the FPR, except in the case of the rare signature when it was noticeably lower (3.7%),
250 presumably due to the asymmetric posterior distribution.

251 We then moved to scenario 1, where we simulated 200 data sets with a difference in mean exposures
252 between the two groups of catalogs. Here, the statistical powers of both HiLDA and the two-stage method
253 are high when detecting the difference in exposures for the 1st and 3rd signatures (Table 4). In contrast,
254 for the 2nd signature, which has the lowest mean mutational exposure, the TPRs of all methods are
255 lower (77.5% - 85.5%). By using the 95% credible interval of posterior distributions, HiLDA is able to
256 distinguish a difference more often than the two-stage method (99.5% vs. 99.0%, 85.5% vs. 77.5%, and
257 91.5% vs. 88.0%). At the same time, using the credible interval resulted in higher TPRs compared to
258 performing a Wald test (85.5% vs. 80.5% for the 2nd signature). In summary, across tests involving these
259 three mutational signatures, HiLDA provides higher statistical power to the TS method with a tendency of
260 better improvement for signatures with lower mutational exposures, i.e., the power difference between
261 HiLDA and the TS method is the highest (8%) for signature 2 with the lowest mean mutational exposures.
262 The improvements in the power to detect the mean exposure difference is presumably due to the fact that
263 HiLDA accounts for the uncertainty in the estimated mutational exposures and provides better model fit
264 of the posterior distributions. All data were simulated in R 3.5.0 using the hierarchical Bayesian mixture
265 model described in the methods section. All replicates reached convergence with an Rhat value less than
266 1.05 for each of the scenarios shown in Tables 2-4.

267 DISCUSSION

268 In this paper, we present a new hierarchical method, HiLDA, that allows the user to simultaneously
269 extract mutational signatures and infer mutational exposures between two different groups of mutational
270 catalogs, e.g., trunk and branch mutations in our example application. Our method is built on the approach
271 of Shiraishi et al. (2015), in which mutational signatures are characterized under the assumption of
272 independence, and it is the first to provide a unified way of testing whether mutational processes differ
273 between groups (here, between early and late stages of tumor growth). As a result, our method allows
274 us to appropriately control the false positive rates while providing higher power by accounting for the
275 accuracy in the estimated mutational exposures.

276 In our analysis of the USC data, which consist of 32 mutational catalogs extracted from tumors from
277 16 CRC patients, our method detected three signatures and indicated a statistically significant difference
278 in mean exposures between groups. Two of the three signatures resemble signatures 7 and 17 found by
279 Shiraishi et al. (2015). But, in addition, we found a novel signature Shiraishi et al. (2015). Signature 7
280 appears significantly more often in trunk mutations, which is consistent with the fact that it has previously
281 been related to aging and trunk mutations have a longer time over which to occur (conceivably over the
282 lifetime of the patient) than do branch mutations (which occur only during tumor growth). The new
283 signature, which occurred more often in low frequency branch mutations, is very similar to a sequencing
284 artifact described by Alexandrov et al. (2018) (cosine similarity = 0.93). We note that, for the USC data,
285 the conclusions obtained from HiLDA were qualitatively the same as those obtained from the TS method.
286 This is likely due to the relatively large effect size here (i.e., the difference of mean exposures between the
287 two groups, divided by the standard errors of same, also known as the signal-to-noise ratio). (Alexandrov
288 et al., 2018).

289 In the simulation study, both HiLDA and the TS approach were applied to datasets consisting of 16
290 tumors simulated under two scenarios to test for between group differences in the mutational exposures of
291 three signature. The results indicated that our unified approach has higher statistical power for detecting
292 differences in exposures for these signatures while controlling the 5% false positive rate. We suspect that
293 the improvement in statistical power is because our unified method explicitly allows for the uncertainty
294 of inferred mutational exposures, while the two-stage method fails to do so since it incorporates only
295 the point estimates of those exposures. In addition, HiLDA provides posterior distributions for each
296 parameter, thereby allowing construction of 95% credible intervals for parameters, and their differences,
297 for example. As expected, this fully parametric approach is then more powerful than nonparametric
298 approaches, which we see particularly when testing for differences in the rarer signatures.

299 We also note that the two-stage approach can become problematic with regards to controlling the type
300 I error rate in particular scenarios, e.g., when the variances of exposures differ widely between the two
301 groups. In our simulation study, we aimed to emulate the USC data, meaning that the exposure variances
302 were quite similar between groups. Consequently, the Wilcoxon rank-sum test, the second-stage of the
303 TS approach, was able to maintain a type I error of 5%. However, we note that the Wilcoxon rank-sum
304 test is sensitive to differences found in either location or scale parameters of the two distributions being
305 tested, i.e., it is sensitive to changes in both the mean and the variance. Therefore, when the variances
306 change between two groups, the Wilcoxon rank-sum test may indicate statistically significant differences
307 in distributions even when the means have not changed, (i.e., due to the difference in shape parameters
308 rather than a difference between location parameters). In contrast, HiLDA explicitly focuses on detecting
309 differences in means, and is robust to effects such as changes in variance. Consequently, when applying
310 the TS method, one should be wary of interpreting significant results as evidence of a "difference in
311 means" when using the TS method (as seems to be common Qin et al. (2018); Meier et al. (2018); Network
312 et al. (2017)). We note that scenarios in which the variance of the estimated exposures differs will be
313 common if the numbers of mutations per tumor varies between the two groups (e.g. when comparing
314 microsatellite instable vs. microsatellite stable colon tumors), leading to an inflated false-positive rate if
315 results from the TS method are interpreted as being evidence of a difference in means. (See Fig. S6. for a
316 specific example of this.) We intend to explore this issue further in a future paper. We also intend to more
317 fully investigate the factors that drive the ability to detect significant difference between groups across a
318 much wider variety of scenarios.

Notation	Description
I	Total number of mutational catalogs (indexed by i)
J_i	Number of observed mutations in i th mutational catalog (indexed by j)
L	Number of features to include. Here, we use the nucleotide substitution, flanking bases and transcription strand (indexed by l)
\mathbf{M}	Vector of the maximum numbers of possible values, (M_1, \dots, M_L) , for each mutation feature, (indexed by M_l), $M_1 = 6$ for nucleotide substitution, $M_2 = 4$ for flanking base, (A, C, G, T), $M_L = 2$ for transcription strand, (+, -)
K	Total number of mutational signatures (indexed by k)
$\mathbf{X}_{i,j}$	Observed mutation characteristic vector, $(x_{i,j,1}, \dots, x_{i,j,L})$, for the j th mutation from the i th mutational catalog (indexed by $x_{i,j,l}$)
$z_{i,j}$	Index of the latent assignment for $\mathbf{X}_{i,j}$, $z_{i,j} \in \{1, \dots, K\}$
$\mathbf{q}_{i,k}$	Probability vector of signature k exposure in mutational catalog i , $(q_{i,1}, \dots, q_{i,K})$, with $\sum_k q_{i,k} = 1$
$\mathbf{f}_{k,l}$	Probability vector of observing any of M_l elements for l th mutation feature, $\mathbf{f}_{k,l} = (f_{k,l,1}, \dots, f_{k,l,M_l})$ with $\sum_{m_l} f_{k,l,m_l} = 1$
\mathbf{F}_k	A tuple of probability vectors with length L , $(\mathbf{f}_{k,1}, \dots, \mathbf{f}_{k,L})$
\mathbf{g}	A vector indicating group membership of the samples. ($g_i \in \{1, 2\}$ for each sample i)
$\boldsymbol{\alpha}$	A tuple of concentration parameters of a Dirichlet distribution with length K , $(\alpha_1, \dots, \alpha_K)$, where the dispersion $\phi = \sum_k \alpha_k$
$\boldsymbol{\mu}$	A tuple of expected values of \mathbf{q} of a Dirichlet distribution with length K , (μ_1, \dots, μ_K) , where $\sum_k \mu_k = 1$.

Table 1. List of notation.

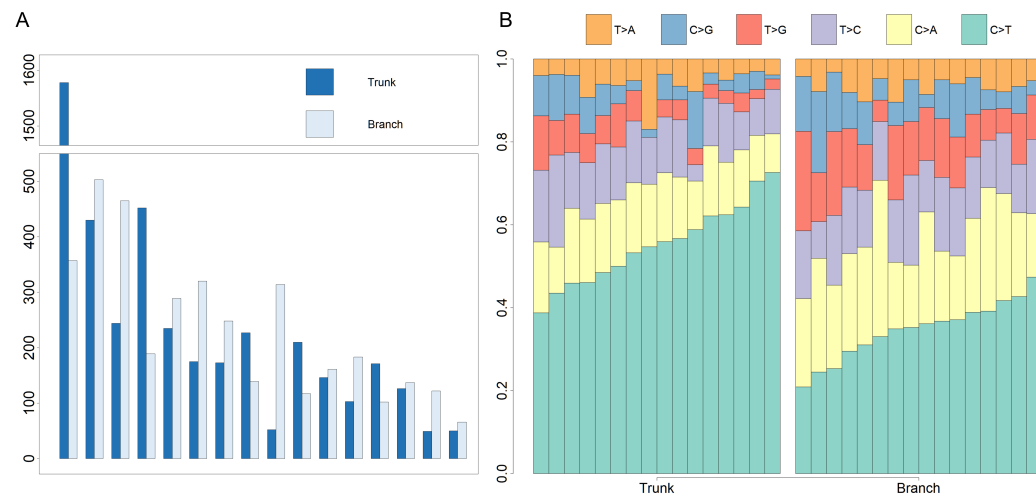


Figure 1. The numbers of somatic mutations in 32 mutational catalogs obtained from 16 colon cancer patients in the USC data and their mutation spectra.

(A) The number of somatic mutations in 16 tumors, each of which contributes 2 mutational catalogs denoted as trunk (dark blue) and branch (light blue).

(B) The percentage bar plot of relative frequencies for six substitution types in the 16 trunk mutational catalogs (left side) and the 16 branch mutational catalogs (right side).

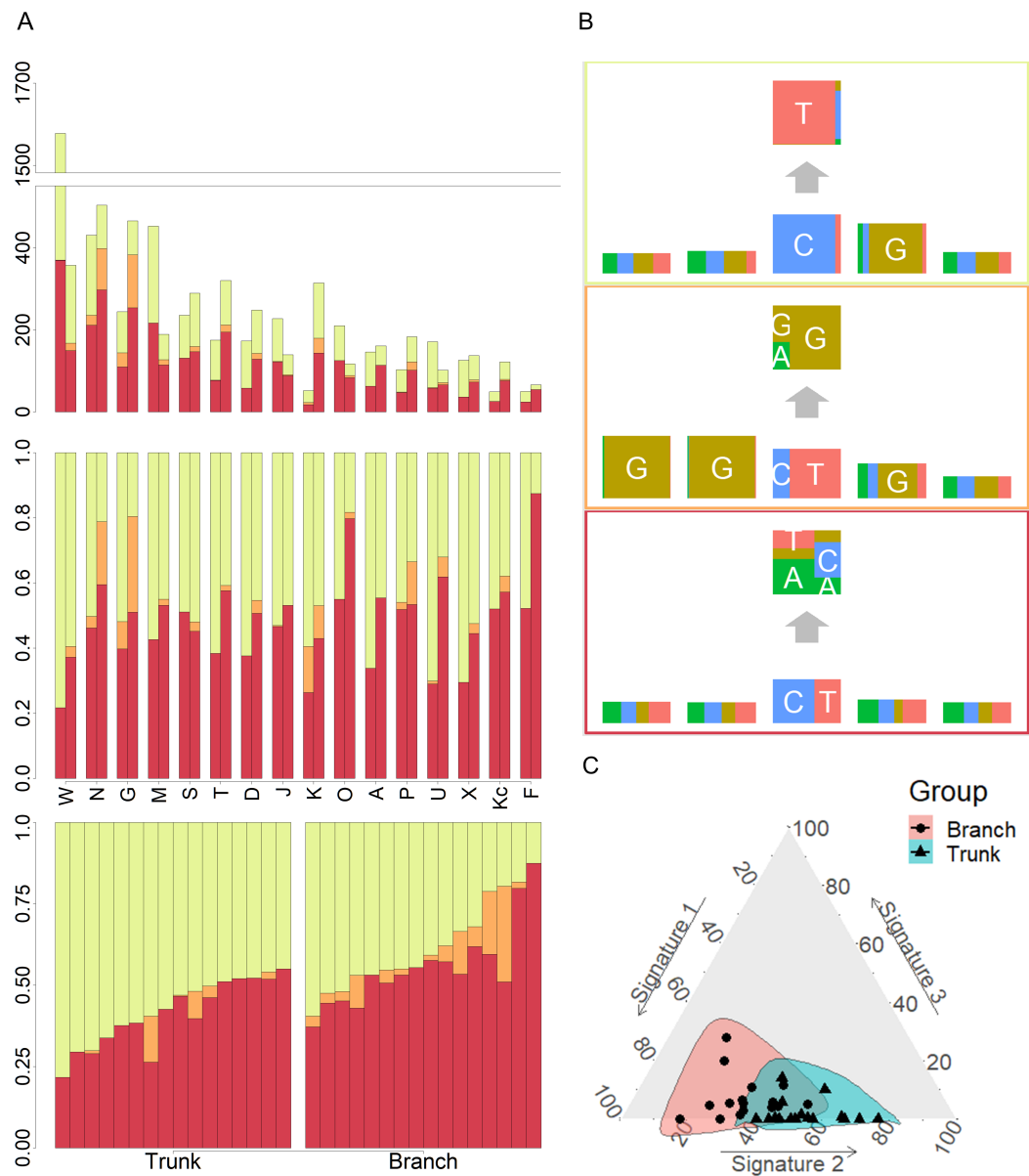


Figure 2. mutational exposures and three mutational signatures from the analysis of 16 trunk mutational catalogs and 16 branch mutational catalogs in the USC data (16 colon cancer patients).

(A) From top to bottom, the three plots represent the somatic mutation counts, the corresponding mutational exposures, and the mutational exposures sorted by group (trunk/branch) and the exposure frequency of the first signature (yellow).

(B) The three mutational signatures with four flanking bases.

(C) The distributions of mutational exposures of the three mutational signatures highlighted by group, where the branch mutational catalogs are highlighted as pink and the trunk ones are highlighted as blue.

	Side A - Side B	HiLDA-CI	HiLDA-Wald	TS-Wilcoxon
Tests ^a	Coef.	[95% C.I.] ^b	p value	p value
Δ_1	0.002	[-0.079, 0.083]	0.9863	0.7804
Δ_2	0.000	[-0.029, 0.029]	0.9875	0.8965
Δ_3	-0.002	[-0.083, 0.086]	0.9608	0.9852
$H_0 : \Delta_1 = \Delta_2 = \Delta_3 = 0$		Bayes Factor _{M₂/M₁} = 0.021		

^a $\Delta_k = \frac{\alpha_k^{(2)}}{\sum_k \alpha_k^{(2)}} - \frac{\alpha_k^{(1)}}{\sum_k \alpha_k^{(1)}}$, the difference in the mean exposure of signature k in group 1 and 2.

^b 95% credible interval from the posterior distribution.

Table 2. Comparing mutational exposures from two sets of mutational catalogs, Side A and Side B, in the USC data.

	Branch - Trunk	HiLDA-CI	HiLDA-Wald	TS-Wilcoxon
Tests^a	Coef.	[95% C.I.]	p value	p value
Δ_1	-0.210	[-0.295, -0.127]	<0.0001	0.0002
Δ_2	0.064	[0.035, 0.099]	0.0001	0.0075
Δ_3	0.146	[0.056, 0.231]	0.0011	<0.0001
$H_0 : \Delta_1 = \Delta_2 = \Delta_3 = 0$		Bayes Factor $_{M_2/M_1} = 1265.0$		

^a $\Delta_k = \frac{\alpha_k^{(2)}}{\sum_k \alpha_k^{(2)}} - \frac{\alpha_k^{(1)}}{\sum_k \alpha_k^{(1)}}$, the difference in the mean exposure of signature k in group 1 and 2.

^b 95% credible interval from the posterior distribution.

Table 3. Comparing mutational exposures in colorectal cancer from two sets of mutational catalogs, trunk and branch, in the USC data.

	Methods	Δ_1	Δ_2	Δ_3
FPRs	HILDA-CI ^a	4.8 %	5.0%	5.1%
	HILDA-Wald ^b	5.1 %	3.7%	5.4%
	TS-Wilcoxon	4.3 %	5.2%	4.3%
TPRs	HILDA-CI	99.5%	85.5%	91.5%
	HILDA-Wald	99.5%	80.5%	92.5%
	TS-Wilcoxon	99.0%	77.5%	88.0%

^a Percentage of 95% credible intervals that exclude zero.

^b Percentage of P -values < 0.05 after applying the Wald test to the posterior distribution.

Table 4. The false positive rates ($n = 1,000$) and true positive rates ($n = 200$) of both the two-stage method and HiLDA when applied to the simulated data.

319 REFERENCES

- 320 Alexandrov, L., Kim, J., Haradhvala, N. J., Huang, M. N., Ng, A. W., Boot, A., Covington, K. R.,
321 Gordenin, D. A., Bergstrom, E., Lopez-Bigas, N., et al. (2018). The repertoire of mutational signatures
322 in human cancer. *bioRxiv*, page 322859.
- 323 Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J., and Stratton, M. R. (2013). Deciphering
324 signatures of mutational processes operative in human cancer. *Cell Reports*, 3(1):246–259.
- 325 Baez-Ortega, A. and Gori, K. (2017). Computational approaches for discovery of mutational signatures
326 in cancer. *Briefings in Bioinformatics*, 20(1):77–88.
- 327 Carlin, B. P. and Chib, S. (1995). Bayesian model choice via markov chain monte carlo methods. *Journal*
328 *of the Royal Statistical Society. Series B (Methodological)*, 57(3):473–484.
- 329 Chang, J., Tan, W., Ling, Z., Xi, R., Shao, M., Chen, M., Luo, Y., Zhao, Y., Liu, Y., Huang, X., et al.
330 (2017). Genomic analysis of oesophageal squamous-cell carcinoma identifies alcohol drinking-related
331 mutation signature and genomic alterations. *Nature Communications*, 8:15290.
- 332 Cibulskis, K., Lawrence, M. S., Carter, S. L., Sivachenko, A., Jaffe, D., Sougnez, C., Gabriel, S.,
333 Meyerson, M., Lander, E. S., and Getz, G. (2013). Sensitive detection of somatic point mutations in
334 impure and heterogeneous cancer samples. *Nature Biotechnology*, 31(3):213–219.
- 335 DePristo, M. A., Banks, E., Poplin, R., Garimella, K. V., Maguire, J. R., Hartl, C., Philippakis, A. A.,
336 Del Angel, G., Rivas, M. A., Hanna, M., et al. (2011). A framework for variation discovery and
337 genotyping using next-generation DNA sequencing data. *Nature Genetics*, 43(5):491–498.
- 338 Gelman, A., Rubin, D. B., et al. (1992). Inference from iterative simulation using multiple sequences.
339 *Statistical Science*, 7(4):457–472.
- 340 Guo, J., Huang, J., Zhou, Y., Zhou, Y., Yu, L., Li, H., Hou, L., Zhu, L., Ge, D., Zeng, Y., et al. (2018).
341 Germline and somatic variations influence the somatic mutational signatures of esophageal squamous
342 cell carcinomas in a chinese population. *BMC Genomics*, 19(1):538.
- 343 Haradhvala, N., Kim, J., Maruvka, Y., Polak, P., Rosebrock, D., Livitz, D., Hess, J., Leshchiner, I.,
344 Kamburov, A., Mouw, K., et al. (2018). Distinct mutational signatures characterize concurrent loss of
345 polymerase proofreading and mismatch repair. *Nature Communications*, 9(1):1746.
- 346 Hillman, R. T., Chisholm, G. B., Lu, K. H., and Futreal, P. A. (2017). Genomic rearrangement signatures
347 and clinical outcomes in high-grade serous ovarian cancer. *JNCI: Journal of the National Cancer*
348 *Institute*, 110(3):265–272.
- 349 Kasuya, E. (2001). Mann-whitney u test when variances are unequal. *Animal Behaviour*, 6(61):1247–1249.
- 350 Letouzé, E., Shinde, J., Renault, V., Couchy, G., Blanc, J.-F., Tubacher, E., Bayard, Q., Bacq, D., Meyer,
351 V., Semhoun, J., et al. (2017). Mutational signatures reveal the dynamic interplay of risk factors and
352 cellular processes during liver tumorigenesis. *Nature Communications*, 8(1):1315.
- 353 Lodewyckx, T., Kim, W., Lee, M. D., Tuerlinckx, F., Kuppens, P., and Wagenmakers, E.-J. (2011). A
354 tutorial on bayes factor estimation with the product space method. *Journal of Mathematical Psychology*,
355 55(5):331–347.
- 356 Mann, H. B. and Whitney, D. R. (1947). On a test of whether one of two random variables is stochastically
357 larger than the other. *The Annals of Mathematical Statistics*, 18(1):50–60.
- 358 Meier, B., Volkova, N. V., Hong, Y., Schofield, P., Campbell, P. J., Gerstung, M., and Gartner, A. (2018).
359 Mutational signatures of DNA mismatch repair deficiency in *C. elegans* and human cancers. *Genome*
360 *Research*, 28(5):666–675.
- 361 Network, C. G. A. R. et al. (2017). Integrated genomic characterization of oesophageal carcinoma. *Nature*,
362 541(7636):169.
- 363 Olivier, M., Bouaoun, L., Villar, S., Robitaille, A., Cahais, V., Heguy, A., Byrnes, G., Le Calvez-Kelm,
364 F., Torres-Mejia, G., Alvarado-Cabrero, I., et al. (2019). Molecular features of premenopausal breast
365 cancers in latin american women: Pilot results from the precama study. *PloS one*, 14(1):e0210372.
- 366 Plummer, M. et al. (2003). Jags: A program for analysis of Bayesian graphical models using Gibbs
367 sampling. In *Proceedings of the 3rd International Workshop on Distributed Statistical Computing*,
368 volume 124, page 125. Vienna, Austria.
- 369 Qin, T., Zhang, Y., Zarins, K. R., Jones, T. R., Virani, S., Peterson, L. A., McHugh, J. B., Chepeha, D.,
370 Wolf, G. T., Rozek, L. S., et al. (2018). Expressed hnscc variants by hpv-status in a well-characterized
371 michigan cohort. *Scientific Reports*, 8(1):11458.
- 372 R Core Team (2017). *R: A Language and Environment for Statistical Computing*. R Foundation for
373 Statistical Computing, Vienna, Austria.

- 374 Ruxton, G. D. (2006). The unequal variance t-test is an underused alternative to student's t-test and the
375 mann–whitney u test. *Behavioral Ecology*, 17(4):688–690.
- 376 Ryser, M. D., Min, B.-H., Siegmund, K. D., and Shibata, D. (2018). Spatial mutation patterns as
377 markers of early colorectal tumor cell mobility. *Proceedings of the National Academy of Sciences*,
378 115(22):5774–5779.
- 379 Saunders, C. T., Wong, W. S., Swamy, S., Becq, J., Murray, L. J., and Cheetham, R. K. (2012). Strelka:
380 accurate somatic small-variant calling from sequenced tumor–normal sample pairs. *Bioinformatics*,
381 28(14):1811–1817.
- 382 Shiraishi, Y., Tremmel, G., Miyano, S., and Stephens, M. (2015). A simple model-based approach to
383 inferring and visualizing cancer mutation signatures. *PLoS Genetics*, 11(12):e1005657.
- 384 Siegmund, K. and Shibata, D. (2016). At least two well-spaced samples are needed to genotype a solid
385 tumor. *BMC Cancer*, 16(1):250.
- 386 Spiegelhalter, D., Thomas, A., Best, N., and Lunn, D. (2003). Winbugs user manual version 1.4.
- 387 Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A., and Sottoriva, A. (2016). Identification of
388 neutral tumor evolution across cancer types. *Nature Genetics*, 48(3):238.