

# 1 Adaptive Somatic Mutations Calls with 2 Deep Learning and Semi-Simulated Data

3 Remi Torracinta<sup>1</sup>, Laurent Mesnard<sup>2</sup>, Susan Levine<sup>4,5</sup>, Rita Shakhovich<sup>3</sup>,  
4 Maureen Hanson<sup>4</sup>, and Fabien Campagne<sup>1,2,\*</sup>

5 <sup>1</sup>The HRH Prince Alwaleed Bin Talal Bin Abdulaziz Alsaud Institute for Computational  
6 Biomedicine, Weill Cornell Medical College, New York, NY, United States of  
7 America; Department of Physiology and Biophysics, Weill Cornell Medical College, New  
8 York, NY, United States of America

9 <sup>2</sup>INSERM UMR1155 et Service des Urgences Néphrologiques et Transplantation Rénale,  
10 APHP, Hôpital Tenon, Paris, France, Sorbonne Universités, UPMC Université Paris 06,  
11 Paris, France

12 <sup>3</sup>Hematology and Oncology Division, Department of Medicine, Weill Cornell Medical  
13 College Department of Medicine, New York, NY 20021 USA

14 <sup>4</sup>Department of Molecular Biology and Genetics, Biotechnology Building, Cornell  
15 University, Ithaca, NY 14853 USA

16 <sup>5</sup>Levine Clinical Practice, 115 E 72nd St, NY 10021 USA

17 <sup>6</sup>Department of Molecular Biology and Genetics, Cornell University. Ithaca, NY 14853  
18 USA

19 \*To whom correspondence should be addressed: [fac2003@campagnelab.org](mailto:fac2003@campagnelab.org)

## 20 ABSTRACT

21 A number of approaches have been developed to call somatic variation in high-throughput sequencing  
22 data. Here, we present an adaptive approach to calling somatic variations. Our approach trains a deep  
23 feed-forward neural network with semi-simulated data. Semi-simulated datasets are constructed by  
24 planting somatic mutations in real datasets where no mutations are expected. Using semi-simulated  
25 data makes it possible to train the models with millions of training examples, a usual requirement for  
26 successfully training deep learning models. We initially focus on calling variations in RNA-Seq data.  
27 We derive semi-simulated datasets from real RNA-Seq data, which offer a good representation of the  
28 data the models will be applied to. We test the models on independent semi-simulated data as well  
29 as pure simulations. On independent semi-simulated data, models achieve an AUC of 0.973. When  
30 tested on semi-simulated exome DNA datasets, we find that the models trained on RNA-Seq data  
31 remain predictive (sens 0.4 & spec 0.9 at cutoff of  $P \geq 0.9$ ), albeit with lower overall performance  
32 (AUC=0.737). Interestingly, while the models generalize across assay, training on RNA-Seq data  
33 lowers the confidence for a group of mutations. Haloplex exome specific training was also performed,  
34 demonstrating that the approach can produce probabilistic models tuned for specific assays and protocols.  
35 We found that the method adapts to the characteristics of experimental protocol. We further illustrate  
36 these points by training a model for a trio somatic experimental design when germline DNA of both  
37 parents is available in addition to data about the individual. These models are distributed with Goby  
38 (<http://goby.campagnelab.org>).

39 Keywords: Deep Learning, Machine Learning, Somatic Variation Caller, Semi-simulated Data

## 40 INTRODUCTION

41 The analysis of high-throughput sequencing data often involves ranking and filtering millions of observed  
42 genomic sites to locate candidates of biological or clinical interest. For instance, many approaches and  
43 associated software tools have been developed to identify sites of somatic variations in the data derived  
44 from the genome of tumors. Methods have been developed for matched DNA samples, where germline  
45 and somatic tumor tissues are both available for an individual Pabinger et al. [2014], Wang et al. [2013].  
46 A few approaches have also been developed for the related, but more challenging problem of identifying

47 somatic variations in matched RNA-Seq data Sheng et al. [2016].

48 Methods to rank and filter candidate sites are often developed using one of two approaches. Early  
49 development following the introduction of a new assay are often ad-hoc and may include the use of  
50 hard filters or other heuristic(s). These methods are widely recognized as being sub-optimal, but are  
51 nevertheless widely applied to new assays. We believe that these methods are used in these instances  
52 because: (1) they do not require a strong knowledge of probabilistic models, (2) they are easy to implement  
53 in software and (3) domain experts can look at results and contribute suggestions to improve the next  
54 iteration of the approach and (4) they enable ranking interesting candidates in the early days of an assay to  
55 demonstrate its biological interest. At these stages, it is more important to identify a few strong candidates  
56 than to obtain optimal sensitivity and specificity across a wide range of datasets.

57 In contrast, probabilistic methods are developed when an assay becomes popular and more data is  
58 produced that requires more sensitive or specific ranking and filtering tools. Developing probabilistic  
59 methods relies on a model of the source of errors in an assay. Developing this model for a new assay can  
60 be a slow process, but once introduced, probabilistic methods frequently outperform ad-hoc approaches  
61 by a wide margin.

62 In this manuscript, we describe a third approach to the ranking and filtering of genomic sites. We  
63 aimed for an approach that (1) would be fast to develop and implement for a new assay (2) would provide  
64 domain experts with the opportunity to contribute to the development and refinement of the approach,  
65 (3) would yield state of the art probabilistic models, (4) can be applied to a wide range of assays and  
66 experimental designs.

67 We initially developed this approach to call somatic variation in matched RNA-Seq samples. RNA-Seq  
68 is a high-throughput sequencing assay that measures gene expression Cloonan et al. [2008], but whose  
69 data can also be used to identify variation in DNA Sheng et al. [2016] (in the parts of the genes that are  
70 expressed at sufficient level to be detectable in a given sample). A few methods have been developed  
71 to call somatic variations in RNA-Seq data, including Sheng et al. [2016] and this task is generally  
72 considered more challenging than calling somatic variations in exome or whole genome data. RNA-Seq  
73 characteristics that make the assay more challenging are (1) base coverage is unequal across the genome  
74 and driven by the expression of the genes encoded at these bases. (2) Splicing complicates genotype  
75 calling by introducing many locations in the genome (exon-exon junctions of expressed genes) where  
76 aligners may mis-align reads to the reference. (3) RNA editing is a set of biological mechanisms that  
77 modifies the sequence of RNA and can be regulated differently in different cell types. Edited bases may  
78 therefore appear as mutations in RNA when DNA is not mutated.

79 The main contribution of this manuscript is to present a new paradigm to develop approaches to  
80 rank and filter genomic sites. In this paradigm, we use semi-simulated datasets to train a probabilistic  
81 model. We evaluate the performance of models trained in this way for calling mutations in RNA-Seq and  
82 exome data, discuss the portability of the model from one assay to the other, and describe the versatility  
83 of the paradigm by training models for an experimental design where DNA from parents of a patient is  
84 available. In contrast to existing approaches, we propose that this new method can be used to quickly  
85 adapt a probabilistic model to specific experimental and analysis protocols.

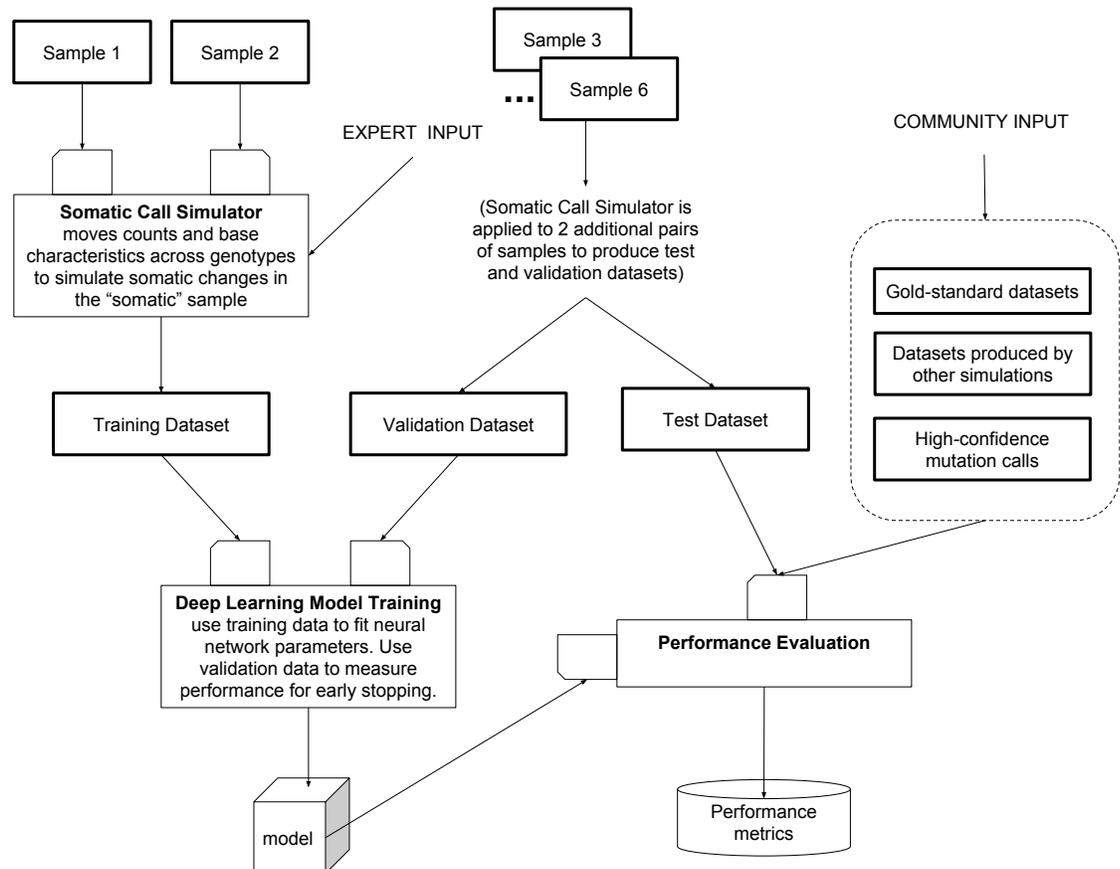
## 86 RESULTS

### 87 A new paradigm to develop probabilistic models

88 The key idea of this study is presented in Figure 1 where we describe how we train probabilistic models  
89 using semi-simulated datasets. This approach does not rely on pre-existing training sets. Instead, we  
90 create semi-simulated datasets by planting signal in real datasets and then train probabilistic models  
91 to recover the signal from the noise found in these datasets. We chose deep learning to implement the  
92 probabilistic models of this approach because deep neural networks (technically, networks are called  
93 “deep” if they have three or more layers) can approximate arbitrary functions and have produced state of  
94 the art performance on a wide range of machine learning problems Angermueller et al. [2016].

### 95 Need for large training sets

96 Training deep learning neural networks requires large training datasets with several orders of magnitude  
97 more examples than parameters in the model. Training state-of-the-art deep learning models often requires  
98 millions of training examples. However, most problems of interest in Bioinformatics have less than a  
99 fraction of a percent of the training size requirements.



**Figure 1. Overview of the approach.** A minimum of two samples are necessary to produce a semi-simulated dataset for training. These samples are chosen from data measured from the same individual such that genotypes should match at the majority of sites measured. One sample is arbitrarily assigned the germline role, and the second sample is assigned the somatic role. These samples are provided to the somatic call simulator (whose design can be informed by expert input), which will produce mutated examples using non-mutated examples provided in the input samples (see Methods). Mutations are always added in the sample labeled “somatic”. This process yields a training dataset. The same process is repeated with independent pairs of input samples (typically from distinct subjects) to yield a validation and a test set. The training set is used to train a feedforward neural network until performance measured on a small part of the validation set starts to decrease (early stopping). Performance of the fit model can be estimated on the validation set as well as on other benchmark datasets contributed by the community.

## 100 **Semi-simulated training sets**

101 To circumvent this problem, we train deep learning models with semi-simulated datasets. Semi-simulated  
102 datasets are constructed by simulating signal into samples where no signal is otherwise expected. For  
103 instance, in this project, we use two RNA-Seq samples from different types of immune cells from the  
104 same subject. Few if any somatic mutations are expected for most sites of the genome in these cells<sup>1</sup>.  
105 We simulate somatic variations in one of the two samples by moving bases and associated features of  
106 the real data to another genotype (see Material and Methods for details). The resulting semi-simulated  
107 dataset should accurately reflect the properties of the real data used as input, and for instance, it is  
108 expected to capture characteristics of the library preparation and sequencing assay used to obtain the data.  
109 Such characteristics include distribution of base quality for errors and somatic variations, distribution of  
110 positions in the reads where variations are observed (i.e., called read index in this study), as well as other  
111 features used to train the models.

## 112 **Training probabilistic models**

113 An overview of this paradigm is shown in Figure 1. Briefly, we train deep neural network models using  
114 large semi-simulated training sets. The training sets are constructed using million of mutation examples  
115 constructed by planting artificial mutations in biological or technical replicates. The replicate samples  
116 are chosen to provide measured data with variability similar to that expected in future samples where the  
117 model will be applied. For instance, to call somatic variations in RNA-Seq data, we trained models using  
118 sorted immune cells from normal controls, where pairs of samples were constructed from a choice of B,  
119 T and NK cells from the same individual. We chose to use biological replicates from different cell types  
120 because somatic variations are often called in a different tissue than the germline tissue and the replicates  
121 need to reflect differences in gene expression between samples.

## 122 **Encoding Data as Features**

123 The samples depicted in Figure 1 consist of reads aligned to a genome. In order to effectively train a  
124 deep learning network, it is necessary to convert alignments to features and labels which are compatible  
125 with back-propagation. In this work, we converted alignment data to features using the process shown in  
126 Figure 2. Briefly, alignments were realigned around indels and lined up by genomic position with the  
127 Goby framework Campagne et al. [2013, 2016b]. This conversion yields summary data about which bases  
128 and indels are observed at a given genomic site, as well as the number of reads that support the base/indel  
129 at the site. These data were then converted to 62 features per site as illustrated in Figure 2 (see Methods  
130 for details about feature mapping).

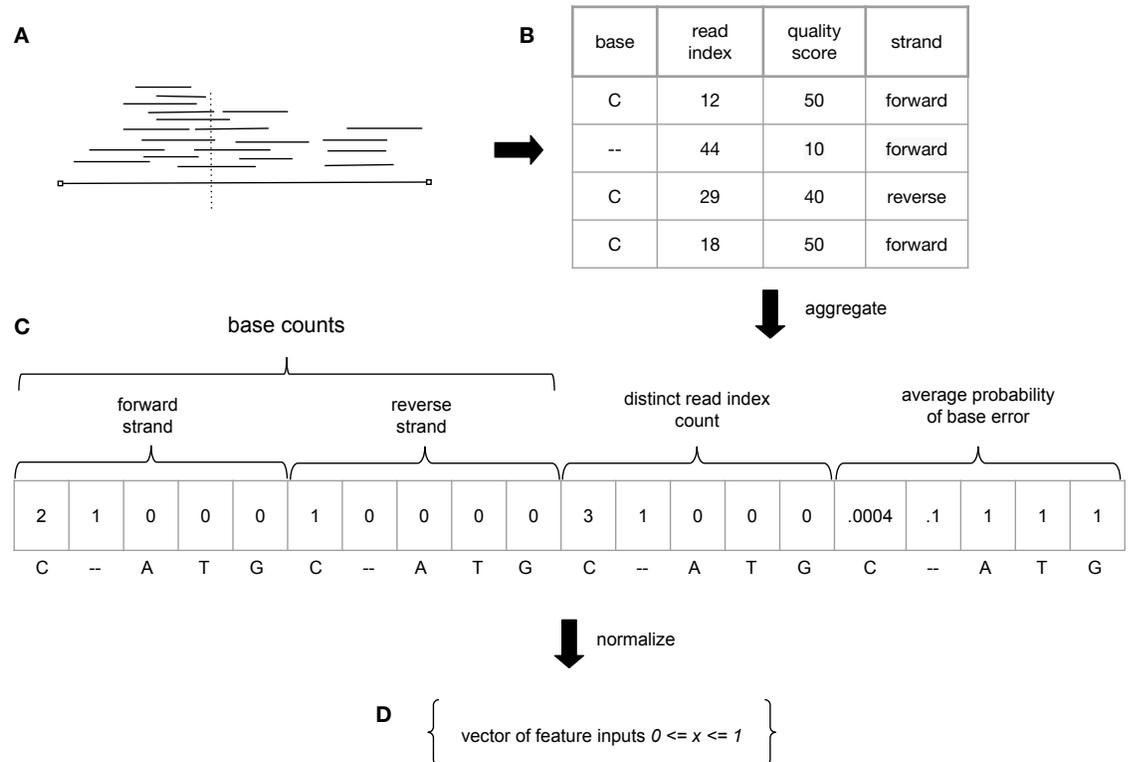
131 The determination of the optimal mapping of read alignment records to features required experimenta-  
132 tion and was performed in an iterative fashion until we found that the score on the validation set (used  
133 for early stopping) did not improve substantially. This iterative model refinement process also included  
134 manual inspection of results obtained on the validation set, to identify any unexpected predictions, and  
135 correct software bugs that can lead to them (see Figure 3). Neural network architecture and neural net  
136 hyper-parameters were also optimized during this process. The final implementation of feature mapping  
137 is modular and makes it possible to plug in or out features derived from an alignment so that different  
138 combinations can be tested and compared.

## 139 **Somatic mutation calling for an RNA-Seq assay**

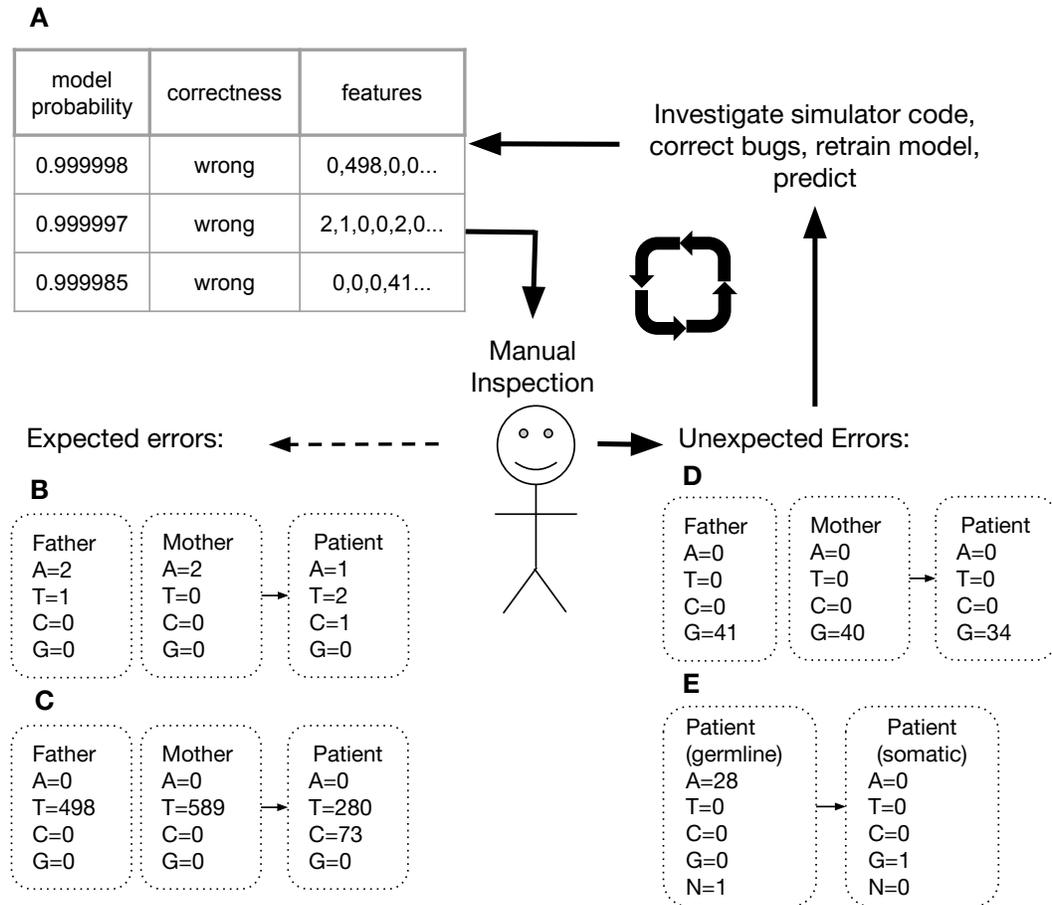
140 We applied the process shown in Figure 1 to develop models to rank and filter somatic mutations in  
141 RNA-Seq data. Briefly, we trained a model with data from two individuals consisting of pairs of control  
142 samples (B vs T, T vs B, NK vs B, B vs NK,...) where mutations were planted to create a semi-simulated  
143 dataset (see Methods for details). The model was trained with early stopping (see Methods). Performance  
144 was measured with the AUC: the area under the Receiver Operator Characteristic (ROC) curve. At the end  
145 of training, performance measured on the validation set was 0.974. A test of the model on an independent  
146 test set gave an AUC of 0.953. Table 1 summarizes the training test and validation performance of the  
147 models developed in this study.

---

<sup>1</sup>Exceptions are the sites that code for B and T cell receptor sites, which undergo somatic mutations in these cell types, but are still a minority among all the sites of a genome.



**Figure 2. Overview of feature encoding.** (A) Aligned reads are processed to obtain data about individual genomic sites (see Material and Methods for details). (B) Illustrates the data collected at each site for observed bases or indels, the genotype of the reference is available, but not shown. (C) While the number of reads at each site can vary dramatically, these data are aggregated into a fixed number of features necessary to train a neural network. The characteristics of each base observed at a site, such as, for instance, average probability of base error (derived from quality scores), number of reads and position in the read supporting the base (read index), are the result of the aggregation. (D) Aggregated features are normalized in the range 0 to 1. This figure shows the steps taken to map one sample to the feature vector (a subset of features is shown). Features derived from additional samples are concatenated.



**Figure 3. Model Refinement Cycle** (A) Example of incorrect predictions on a dataset. The table is ordered by model probability (that the site is a somatic mutation) and filtered to show only false positive predictions. Manual inspection of such a table helps identify the most extreme errors made by a model. (B) A low signal to noise ratio, as in this case, is an expected source of false positives and may indicate that a model needs further training. (C) This negative example is indistinguishable from a positive example (it looks just like the simulator planted a mutation, when this was not the case). Such errors are possible if the genotype of germline samples differ at some sites (D) This negative example was wrongly classified as a mutation. Upon further investigation of the simulator code, we identified a software bug. During the feature encoding phase, the order of bases was not consistent from sample to sample. Fixing this bug resolved such errors. (E) In a first attempt at simulation, we trained with datasets restricted to have at least 10 counts across the samples at a site. This prevented the model from learning that such sites are more difficult to predict and resulted in over-estimates of model probability. Removing filters from the construction of the training set (to include sites irrespective of their coverage) resolved this issue.

Model	Training Set	Validation AUC	Test Set	Test AUC
RNA-Seq	RNA-Seq	0.974	RNA-Seq	0.953
RNA-Seq	RNA-Seq	0.974	Pair Exome	0.737
Pair Exome	Pair Exome	0.996	Pair Exome	0.994
Trio Exome	Trio Exome	0.990	Trio Exome	0.993

**Table 1.** Performance of Models

#### 148 **Somatic mutation calling for an exome assay**

149 An important question is whether the modeling infrastructure developed and optimized for the RNA-Seq  
150 assay can be reused to train models for a different assay.

151 To address this question, we trained a new model with exome data obtained with the HaloPlex exome  
152 assay (Agilent, see <http://www.genomics.agilent.com/article.jsp?pageId=3081>). This HaloPlex assay uses  
153 enzymatic cleavage of DNA which results in very different alignment profiles than RNA-Seq protocols  
154 (which usually employ sonication to break cDNAs into fragments before adapter ligation). HaloPlex  
155 reads stack sharply at the positions where the mix of restriction enzymes cut DNA. To train the model  
156 with exome data, we used two germline samples from the same individual. These samples are expected  
157 to display germline genotypes across the entire exome. We produced semi-simulated training sets and  
158 trained an HaloPlex exome model. At the end of training, performance measured on the validation set was  
159 0.996. A test of the model on an independent test set gave an AUC of 0.994 (see Table 1). The increase in  
160 performance of the model was expected because an exome assay is specifically optimized to produce data  
161 suitable to call genotypes, whereas RNASeq assays are optimized to measure gene expression. This result  
162 suggests that the features developed for the RNA-Seq data are transferable to a different assay when a  
163 model is retrained with a dataset for the new assay.

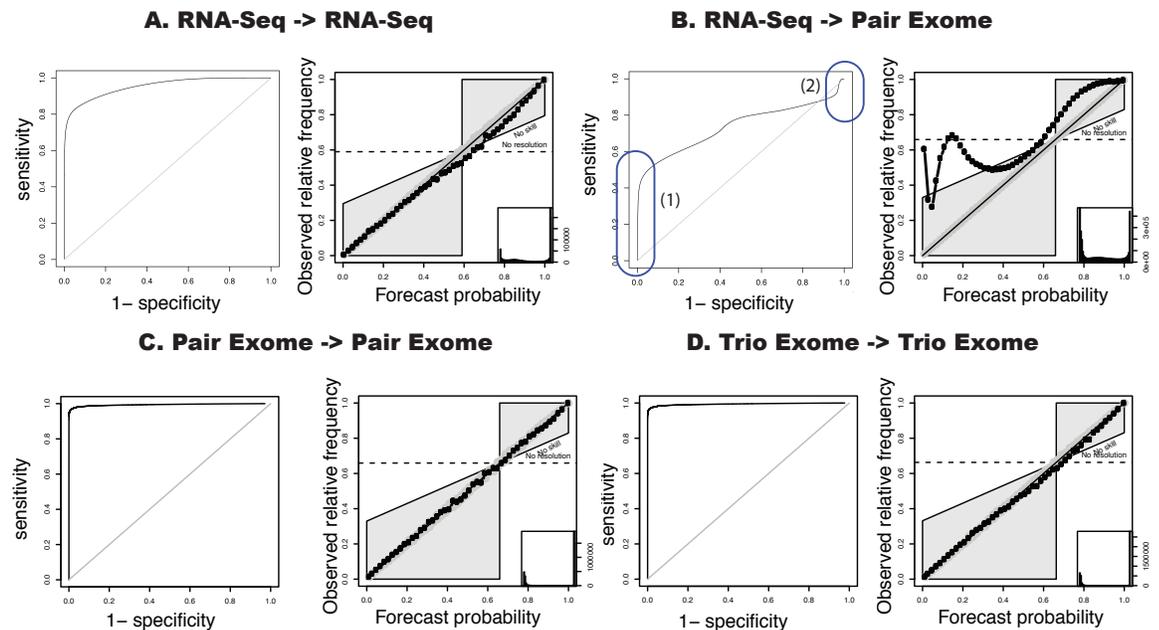
#### 164 **Somatic mutation calling for exome and trio design**

165 An additional question is whether the semi-simulation approach can be adapted to support different  
166 experimental designs.

167 To address this question, we developed models for a trio design. Assume somatic calls are needed  
168 for a patient for whom germline DNA is not readily available (because blood is suspected to contain  
169 cells with the somatic variations and contamination of other tissues by somatic cells cannot be ruled out).  
170 Assuming data from DNA from both parents of the patient are available, it should be possible to call  
171 somatic variations in the patient sample. We call this experimental design “trio somatic”, in contrast to  
172 the “pair somatic” design discussed in the previous section. We combined feature mappers developed for  
173 the pair somatic problem such that features that depend on a single sample are calculated successively for  
174 the father, mother and patient, and features that depend on two samples are calculated for father/patient  
175 and mother/patient pairs. All features were concatenated. As shown in Table 1, training this model  
176 produced similar performance to that trained for the pair exome design. This result indicates that the  
177 semi-simulation approach can train models that learn the mendelian inheritance rules necessary to call  
178 somatic variations in trio designs.

#### 179 **Adaptive Models**

180 The performance of the models developed in this study is summarized in Figure 4. Each panel shows the  
181 received operating curve (ROC) as well as a reliability diagrams Niculescu-Mizil and Caruana [2005].  
182 The reliability diagrams compare the expected true positive rate given by the output model probability, to  
183 that observed in an independent test set(Niculescu-Mizil and Caruana [2005]). In each panel (A, C and D)  
184 where we train the model on the same experimental protocol, we find almost optimal model reliability.  
185 In Panel B, we test a model trained on RNA-Seq data on the paired exome data. The purpose of this  
186 comparison is to determine how transferable a model trained in one assay is to another assay. We find  
187 that the model remained predictive for about 40% of the exome sites, but has lower performance on the  
188 exome data. The reliability diagram also confirms that the model performs less reliably across assays.  
189 Taken together these data suggest that large performance gains can be obtained by training a model to the  
190 specific assay that it will be used for, adapting the model to the assay. While we have not measured this  
191 effect at this time, we expect that smaller performance advantages may result from adapting a model to



**Figure 4. Model Performance** We characterize performance with the Receiver Operating Curves (ROC, shown on the left of each panel) and Reliability Diagrams (RC, shown on the right of each panel). (A) The model trained on RNA-Seq data performs well on RNA-Seq data and has strong reliability. (B) The same model trained on Pair Exome data predicts accurately a subset of high-confidence sites (1), with high predicted probability of mutation, but has degraded performance on other sites (2). This is confirmed by the RC, which shows sub-optimal reliability.

192 specific data analysis protocols.

### 193 **Software implementation**

#### 194 **Using trained models**

195 Models trained in this study have been integrated in release 3.0+ of the Goby somatic caller (distributed  
196 at <http://github.com/CampagneLaboratory/goby3>). Goby3 supports alignments in the  
197 Goby or BAM formats. A parameter is used to specify the path to a model to estimate probabilities of  
198 somatic variation.

#### 199 **Training new models**

200 Models for new assays can be trained by constructing a semi-simulated dataset using data obtained with the  
201 specific assay or analysis protocol. The samples used to train the model should contain no somatic variation  
202 (or a small expected number of variations when completely germline samples cannot be obtained). Detailed  
203 steps are documented with the software (see [https://github.com/CampagneLaboratory/  
204 variationanalysis](https://github.com/CampagneLaboratory/variationanalysis)), but briefly, the semi-simulated dataset is produced by converting the samples  
205 to a raw dataset using the Goby SEQUENCE\_BASE\_INFORMATION output format. The file is mutated  
206 with the simulator corresponding to the experimental design (pair or trio), randomized, and split into  
207 training, validation and test sets. The training of the model uses the training and validation sets. Final  
208 model performance can be measured on the test set to check that the model generalizes.

## 209 **DISCUSSION**

### 210 **Training probabilistic methods**

211 We have presented a new approach to develop probabilistic models for calling somatic variations in high-  
212 throughput sequencing data. In contrast to previous studies which have manually designed probabilistic  
213 models, we show that it is possible to train probabilistic models using semi-simulated datasets and deep  
214 learning methods and that such models adapt to the characteristics of the data produced by specific assays  
215 and analysis protocols.

## 216 **Additional validations are needed**

217 Additional validations will be needed to firmly establish that the models trained in this way are predictive  
218 on real datasets. Because of regulatory limitations on data sharing of genotype data, which have hampered  
219 our ability to obtain suitable validation data, we chose to distribute the software and the models to make it  
220 possible for others to conduct independent evaluations on private datasets. We hope to be able to test the  
221 models on the ICGC GoldSet Alioto et al. [2015] (request initiated with ICGC in July 2016, approved  
222 Sept 2016, we have yet to obtain access to the data files for the ICGC GoldSet which cannot be located  
223 through the ICGC portal. An email to the ICGC support mailing list to request access through the ENA,  
224 as per ICGC instructions on the ENA web site has remained unanswered as of end of Sept 2016).

## 225 **Training for new Assays and Experimental Designs**

226 Should our approach perform competitively on real datasets, one of its key advantage will be that training  
227 can be performed for arbitrary combinations of new assays, experimental designs, and analysis protocol.

228 For instance, training for a new assay only requires a dataset where no somatic variations are expected  
229 (such as biological or technical replicates from the same individual). An assay-specific model can be  
230 trained following simple documented steps using the software that we have developed for this study.  
231 Training for a different analysis protocol (e.g., combination of read pre-processing and aligner) can be  
232 performed similarly.

233 Training for a new experimental design requires additional feature engineering, which must be  
234 implemented in new software. In our experience, developing and testing new feature mappers can be done  
235 in a couple weeks by a scientist familiar with the software when adequate computational resources are  
236 available.

## 237 **Semi-simulated datasets**

238 Many previous studies have taken advantage of machine learning approaches to train probabilistic methods  
239 using annotated real datasets. This work differs in that we create the training sets using real datasets  
240 where no difference exists, and plant signal artificially in the real data background. If no datasets can  
241 be found with no expected difference, training could be also be performed with a dataset where only a  
242 minority of the sites are mutated. This would result in under-estimates of the probability of mutation at a  
243 site, but this bias should be small if the proportion of mutated sites in the sample is low with respect to the  
244 number of planted mutations introduced in the sample.

## 245 **Of the expected mutation rate**

246 We showed that our approach has strong reliability when applied on the same type of assay as that used in  
247 the training set. The model outputs the probability that a site is mutated, given the prior distribution of  
248 mutated sites in the training set and the reliability diagrams indicate that the forecast probability produced  
249 by the model is close to optimal (i.e., lying very close to the diagonal on the reliability diagram). When  
250 the proportion of mutated sites in a dataset differs from that used in the training set, as is often the case in  
251 practice, it will be necessary to apply Bayes Theorem to adjust for the difference in priors. Doing so will  
252 require an estimate of the rate of mutation in the dataset where predictions must be made. Importantly,  
253 such adjustments will not change rank order, and for this reason we expect that the model probability  
254 output by the software is suitable for ranking somatic mutation in new samples.

## 255 **MATERIAL AND METHODS**

### 256 **Subject characteristics and recruitment**

257 Data from four subjects was used for developing the RNA-Seq somatic caller. RNA was obtained from  
258 sorted cells for control subjects who participated to the CFS/ME study. Any minor subject(s) (<18  
259 years old in New York State) who did not provide written acknowledgment of parental permission to  
260 participate in the study were excluded. The study was reviewed and approved by the Weill Cornell  
261 Medical College Institutional Review Board (protocol 1302013563. "Immune cell gene expression and  
262 predictive models in CFS"), and patients and controls gave written informed consent after the study  
263 protocol was fully explained. All consented to blood draw and to the availability of the stored samples for  
264 additional bioassays and analyses.

265 Data from 12 subjects (three subjects for paired exome caller and nine subjects for trio design:  
266 three subjects and their parents) was used to develop the trio somatic caller. The study involving the

267 exome and trio subjects was approved by the Comité de Protection des Personnes (CPP), Ile de France 5,  
268 (05/12/2012).

### 269 **RNA-Seq alignments**

270 RNA-Seq reads were aligned to the 1000 Genome Project human reference sequence (corresponding to  
271 hg19) using STARR Dobin et al. integrated with GobyWeb Dorff et al. [2013].

### 272 **Exome and Trio alignments**

273 Exome and Trio reads were aligned to the 1000 Genome Project human reference sequence as previously  
274 described Mesnard et al. [2016].

### 275 **Training, Validation and Test Datasets**

276 The RNA-seq training, validation and test sets were created by serializing alignments to the .sbi/.sbip  
277 format with the Goby framework. The .sbi/.sbip format is a protocol buffer format developed with  
278 methods described in Campagne et al. [2013], which serializes information about aligned bases for one or  
279 more samples. The record of information in the .sbi/.sbip format is a single genomic site. For the somatic  
280 models, each record stored information about a sequenced position in two different samples (germline  
281 or somatic), or three samples (mother, father, and somatic). The protobuf schema for the .sbi/.sbip  
282 format is distributed in the Goby 3.0 repository ([https://github.com/CampagneLaboratory/  
283 goby3/goby-distribution/protobuf/BaseInformationRecords.proto](https://github.com/CampagneLaboratory/goby3/goby-distribution/protobuf/BaseInformationRecords.proto)) Campagne  
284 et al. [2016b]. The Goby 3.0 discover-sequence-variation mode was used to serialize Goby alignments to  
285 training, validation and test datasets using the SEQUENCE\_BASE\_INFORMATION format (as described  
286 at <https://github.com/CampagneLaboratory/variationanalysis>).

### 287 **Mutation Simulator**

288 The simulator is responsible for constructing millions of both positive and negative examples to train the  
289 neural network, using data records about individual genomic sites. Data records are produced by Goby  
290 and contain most of the sequencing data collected from the two samples aligned at a given genomic site  
291 (single position). The construction of mutated examples proceeds by scanning data records one at a time.  
292 First, the simulator determines if a record follows the count distribution expected of a diploid genome  
293 (canonical sites). Following this determination, only the canonical sites are mutated. The non-canonical  
294 sites are written as is with the label non-mutated. The criteria to identify canonical sites are as follows:

- 295 • the data indicate that at most two alleles are observed in either sample. This determination is made  
296 when the counts from more than 2 genotypes, with genotypes ordered in decreasing count order,  
297 have to be summed to reach 90% of all base counts for the site.
- 298 • the alleles identified by the 2 genotype rule match across samples at the site.

299 In this study, 18% of RNA-seq and 6% of Haloplex records were found to be non-canonical. Canonical  
300 sites are used by the simulator to create two mutated versions of the data record. The one unmutated and  
301 two mutated versions are added to the training set. Making a mutated version of a record consisted of a  
302 few steps. We used a simple heuristic to determine the original genotype, and picked a random “source”  
303 base from the two alleles in the genotype. We generated a frequency of mutation,  $f$ , from the uniform  
304 distribution between 0 and 0.5 (heterozygote site) or 1 (homozygote site), and chose a random “target”  
305 genotype from any of the genotypes not marked as source. We then subtracted a proportion  $f$  of bases  
306 from the source genotype, and added them to the target genotype. Each base retained the same features it  
307 was associated with in the source base, namely its read index location, quality score, and forward/reverse  
308 read direction.

### 309 **RNA-Seq Datasets**

310 We used B, NK, and T-cell samples from two control subjects (CFS/ME study). For each position, a  
311 different record represented a different permutation of two cell-types for a subject, so we had 12 records  
312 for each position. We introduced mutated records with the mutation simulator, as described previously,  
313 and then randomized the order of the records within the training set.

314 The validation and test sets were created in the same way, using only NK and B cells data from one  
315 subject for each set.

	# Training Examples	# Validation Examples	# Test examples	# Features
<b>RNA-seq</b>	8159893	721012	666090	62
<b>Pair Exome</b>	2776689	3648105	1092657	62
<b>Trio Exome</b>	4487139	3907871	1576158	90

**Table 2. Dataset Characteristics**

### 316 **Pair Exome Datasets**

317 For pair exome datasets, we only created records of positions for one sample pairing (DNA extracted from  
318 the subject's blood and skin). The training, validation, and test sets each drew data from an independent  
319 subject. Mutated records were introduced in training, validation and test sets. The records of each set  
320 were shuffled with the Randomizer2 tool provided in Torracinta and Campagne [2016].

### 321 **Trio Exome Datasets**

322 For trio exome datasets, we created records of positions for one sample trio (DNA from father's blood,  
323 mother's blood, and subject's blood). The training, validation, and test sets drew data from independent  
324 trios. Mutated records were introduced in training, validation and test sets. The records of each set were  
325 shuffled.

### 326 **Neural Network Architecture**

#### 327 **Feature Mappers**

328 Feature mappers convert alignments about one or more samples into a fixed set of features suitable  
329 for training with neural networks. Regardless of the number of reads at a genomic position, mappers  
330 needed to produce a fixed-length output so that these outputs could be concatenated consistently into  
331 a fixed-length input vector. At each position, a mappers were used to generate the number of reads  
332 supporting each genotype (counts), the number of distinct locations in the read that support the genotype  
333 (distinct read indices), the average error probability of a base being an error (derived from base quality  
334 scores at positions that do not match the reference), and the difference in genotype proportion between  
335 the first sample (always germline) and the second sample (possibly containing a somatic mutation). For  
336 every given position, mappers which produced features about just one sample were reapplied to each  
337 sample in the assay with the results concatenated. The last mapper, which generated a comparative feature  
338 between two samples, was applied once for pair design (germline/somatic), and twice for trio design  
339 (father/child,mother/child). Mappers are implemented in the variationanalysis project available at GitHub  
340 <https://github.com/CampagneLaboratory/variationanalysis>.

### 341 **Model Architecture**

342 Models were developed with the DeepLearning4J (DL4J) framework (<http://deeplearning4j.org/>), ver-  
343 sion 0.5.0. DL4J was selected because it is a Java framework and the models it produces can be  
344 integrated with the Goby framework more easily than frameworks in other languages. Models were  
345 formulated as 5 fully connected layers with RELU activation and a fully connected layer with soft-max  
346 activation. The first layer contained 5 times the number of input features. Inner layers contained 0.65  
347 the number of neurons in the preceding layer. The exact model architecture used is encoded in the  
348 class called `org.campagnelab.dl.varanalysis.learning.architecture.SixDense`  
349 `LayersNarrower2` distributed in the variationanalysis project Torracinta and Campagne [2016].

### 350 **Training Procedure**

351 Stochastic gradient descent optimization was used, with an ADAGRAD optimizer (called updater in the  
352 DL4J framework). The model trained on mini-batches of 600 examples. Regularization was not used.

353 Finding that it did not introduce training instability, the learning rate was set to the starting value  
354 of 1, with SCORE decay (the learning rate was decreased when the loss on the validation set stopped  
355 decreasing).

356 An early stopping condition of 10 epochs interrupted training early if no performance improvement  
357 was observed in the validation set after 10 training epochs. Models were trained with early stopping by  
358 measuring performance on 100,000 examples from the validation set.

## 359 Performance measurements

360 AUC measures were calculated with the `org.campagnelab.dl.varanalysis.stats.AreaUnderTheROCCurve`  
361 class, provided in the in the variationanalysis project Torracinta and Campagne [2016]. This class  
362 implements the naive  $O(n^2)$  calculation of the area under the ROC curve and directly calculates the  
363 probability of correctly classifying an observation. This class was adapted from the AUC calculator  
364 implemented in the BDVal project, which was validated in the MAQC-II project Consortium et al. [2010].

365 Received Operating Curves were produced with the AUC R package, integrated in MetaR Campagne  
366 et al. [2016a]. Reliability diagrams were constructed with the SpecsVerification R package, integrated in  
367 MetaR Campagne et al. [2016a].

## 368 AUTHOR CONTRIBUTIONS

369 RT and FC wrote the variation analysis programs as well as parts of Goby3 necessary to create datasets.  
370 LM provided exome and trio data and developed detailed protocols for blood collection and processing in  
371 the CFS/ME study. SL helped obtain the four control samples of the CFS/ME project used to train and  
372 validate models in this study. RS developed detailed protocols for B cell, T cell and NK cell sorting in the  
373 CFS/ME study and supervised blood extraction, cell sorting and RNA extraction. MH and FC supervised  
374 the human subject components of the CFS/ME study. FC designed the computational approaches.

## 375 ACKNOWLEDGMENTS

376 We thank Manuele Simi for technical assistance with the code needed for this project.

## 377 FUNDING

378 This investigation was supported by the National Institutes of Health NIAID award 5R01AI107762 to  
379 Fabien Campagne and Maureen Hanson. This investigation was also supported by the STARR cancer  
380 consortium award I9-A9-084 to Samie Jaffrey, Jedd Wolchok and Fabien Campagne.

## 381 REFERENCES

- 382 Tyler S Alioto, Ivo Buchhalter, Sophia Derdak, Barbara Hutter, Matthew D Eldridge, Eivind Hovig,  
383 Lawrence E Heisler, Timothy A Beck, Jared T Simpson, Laurie Tonon, et al. A comprehensive assess-  
384 ment of somatic mutation detection in cancer using whole-genome sequencing. *Nature communications*,  
385 6, 2015.
- 386 Christof Angermueller, Tanel Pärnamaa, Leopold Parts, and Stegle Oliver. Deep Learning for Computa-  
387 tional Biology. *Molecular Systems Biology*, (12):878, 2016.
- 388 Fabien Campagne, Kevin C. Dorff, Nyasha Chambwe, James T. Robinson, and Jill P. Mesirov. Com-  
389 pression of Structured High-Throughput Sequencing Data. *PLoS ONE*, 8(11):e79871, nov 2013.  
390 ISSN 1932-6203. doi: 10.1371/journal.pone.0079871. URL [http://dx.plos.org/10.1371/](http://dx.plos.org/10.1371/journal.pone.0079871)  
391 [journal.pone.0079871](http://dx.plos.org/10.1371/journal.pone.0079871).
- 392 Fabien Campagne, William ER Digan, and Manuele Simi. Metar: simple, high-level languages for data  
393 analysis with the r ecosystem. *bioRxiv*, page 030254, 2016a.
- 394 Fabien Campagne, Remi Torracinta, and Manuele Simi. Goby 3.0.0 software release, 2016b. URL  
395 <https://doi.org/10.5281/zenodo.159024>.
- 396 Nicole Cloonan, Alistair R R Forrest, Gabriel Kolle, Brooke B a Gardiner, Geoffrey J Faulkner, Mellissa K  
397 Brown, Darrin F Taylor, Anita L Steptoe, Shivangi Wani, Graeme Bethel, Alan J Robertson, Andrew C  
398 Perkins, Stephen J Bruce, Clarence C Lee, Swati S Ranade, Heather E Peckham, Jonathan M Manning,  
399 Kevin J McKernan, and Sean M Grimmond. Stem cell transcriptome profiling via massive-scale mRNA  
400 sequencing. *Nature methods*, 5(7):613–9, 2008. ISSN 1548-7105. doi: 10.1038/nmeth.1223. URL  
401 <http://www.ncbi.nlm.nih.gov/pubmed/18516046>.
- 402 MAQC Consortium et al. The microarray quality control (maq)-ii study of common practices for the  
403 development and validation of microarray-based predictive models. *Nature biotechnology*, 28(8):  
404 827–838, 2010.
- 405 A Dobin, C A Davis, F Schlesinger, J Drenkow, C Zaleski, S Jha, P Batut, M Chaisson, and T R Gingeras.  
406 doi: bts635[pii]10.1093/bioinformatics/bts635.

- 407 Kevin C. Dorff, Nyasha Chambwe, Zachary Zeno, Manuele Simi, Rita Shaknovich, and Fabien Campagne. GobyWeb: Simplified Management and Analysis of Gene Expression and DNA Methylation Sequencing Data. *PLoS ONE*, 8(7):e69666, 2013. ISSN 1932-6203. doi: 10.1371/journal.pone.0069666. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3720652&tool=pmcentrez&rendertype=abstract>.
- 412 Laurent Mesnard, Thangamani Muthukumar, Maren Burbach, and Carol Li. Exome Sequencing and Prediction of Long- Term Kidney Allograft Function. pages 1–15, 2016. doi: 10.1371/journal.pcbi.1005088.
- 415 Alexandru Niculescu-Mizil and Rich Caruana. Predicting good probabilities with supervised learning. In *Proceedings of the 22nd international conference on Machine learning*, pages 625–632. ACM, 2005.
- 417 Stephan Pabinger, Andreas Dander, Maria Fischer, Rene Snajder, Michael Sperk, Mirjana Efre-mova, Birgit Krabichler, Michael R Speicher, Johannes Zschocke, and Zlatko Trajanoski. A survey of tools for variant analysis of next-generation genome sequencing data. *Briefings in bioinformatics*, 15(2):256–78, mar 2014. ISSN 1477-4054. doi: 10.1093/bib/bbs086. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3956068&tool=pmcentrez&rendertype=abstract>.
- 423 Quanhu Sheng, Shilin Zhao, Chung I. Li, Yu Shyr, and Yan Guo. Practicability of detecting somatic point mutation from RNA high throughput sequencing data. *Genomics*, 107(5):163–169, 2016. ISSN 10898646. doi: 10.1016/j.ygeno.2016.03.006.
- 426 Remi Torracinta and Fabien Campagne. Variationanalysis 1.0.2 software release, October 2016. URL <https://doi.org/10.5281/zenodo.159203>.
- 428 Qingguo Wang, Peilin Jia, Fei Li, Haiquan Chen, Hongbin Ji, Donald Hucks, Kimberly Brown Dahlman, William Pao, and Zhongming Zhao. Detecting somatic point mutations in cancer genome sequencing data: a comparison of mutation callers. *Genome medicine*, 5(10):91, 2013. ISSN 1756-994X. doi: 10.1186/gm495. URL [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3971343&tool=pmcentrez&rendertype=abstract&backslash\\$nhhttp://www.ncbi.nlm.nih.gov/pubmed/24112718&backslash\\$nhhttp://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3971343](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3971343&tool=pmcentrez&rendertype=abstract&backslash$nhhttp://www.ncbi.nlm.nih.gov/pubmed/24112718&backslash$nhhttp://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3971343).