

1 **OPENANNO: annotating genomic regions with chromatin**
2 **accessibility**

3

4 Shengquan Chen^{1,a}, Yong Wang^{2,3,*b}, Rui Jiang^{1,*c}

5

6 ¹ *MOE Key Laboratory of Bioinformatics; Bioinformatics Division and Center for*
7 *Synthetic and Systems Biology; Beijing National Research Center for Information*
8 *Science and Technology; Department of Automation, Tsinghua University, Beijing*
9 *100084, China*

10 ² *CEMS, NCMIS, MDIS, Academy of Mathematics and Systems Science, Chinese*
11 *Academy of Sciences, Beijing 100190, China*

12 ³ *Center for Excellence in Animal Evolution and Genetics, Chinese Academy of*
13 *Sciences, Kunming 650223, China.*

14

15 * Corresponding authors.

16

17 **Abstract**

18 Chromatin accessibility, as a powerful marker of active DNA regulatory elements,
19 provides rich information to understand the regulatory mechanism. The revolution in
20 high-throughput methods has accumulated massive chromatin accessibility profiles in
21 public repositories as a valuable resource for machine learning and integrative studies.
22 Nevertheless, utilization of these data is often hampered by the cumbersome and time-
23 consuming collection, processing, and annotation of the chromatin accessibility
24 information. Motivated by the above understanding, we developed a web server, named
25 OPENANNO, to **annotate** the **openness** of genomic regions across diverse cell lines,
26 tissues, and systems. The annotation is based on 871 DNase-seq experiments across
27 199 cell lines, 48 tissues, and 11 systems from ENCODE, and openness values
28 rigorously defined by four statistical strategies. Particularly, we designed a parallel
29 program to allow efficient annotation and visualization of the openness of a vast amount
30 of genomic regions. OPENANNO will help users extract and download formulated data
31 in a batch follow-up analysis. Besides, we illustrate the valuable information provided
32 by OPENANNO using an enhancer of blood vessels from VISTA Enhancer Browser
33 as an example. Furthermore, we demonstrate three applications of OPENANNO in
34 regulatory mechanism and association studies. We believe that OPENANNO will serve
35 as a comprehensive and user-friendly web server to facilitate methodology
36 development and biological insights discovery, specifically to explore the biological
37 questions and model the regulatory landscape of genome. OPENANNO is freely
38 available at <http://bioinfo.au.tsinghua.edu.cn/openness/anno> or
39 <http://159.226.47.242:65424/openness/anno/>.

40

41 **KEYWORDS:** Chromatin accessibility; Openness; Annotation; Visualization; Web
42 server

43

44 **Introduction**

45 The era of the personal genome is arriving with the widespread sequencing technologies
46 and the ultimate promise for precision medicine. However, it remains distant in
47 interpreting the context of variations in non-coding DNA sequence associated with
48 disease and other phenotypes, deciphering their biological functions in gene regulation,
49 and further understanding disease mechanism and dynamic response to treatment [1, 2].
50 Chromatin accessibility is a measure of the ability of nuclear macromolecules to
51 physically contact DNA [3], and plays the role of a powerful marker of active regulatory
52 genomic regions, which have a wide range of effects on the transcription, DNA repair,
53 recombination, and replication [4, 5]. The recent revolution in high-throughput,
54 genome-wide methods invented several biological assays for extracting open chromatin,
55 such as DNase-seq (deoxyribonuclease), FAIRE-seq (formaldehyde-assisted isolation
56 of regulatory elements), ATAC-seq (assay for transposase-accessible chromatin) and
57 MNase-seq (micrococcal nuclease), and thus open a new door for us to make extensive
58 use of chromatin accessibility [6-11]. For example, accessible genomic regions are
59 regarded as the primary positions of regulatory elements [12], and thus provide a great
60 opportunity to study transcription factor binding sites, DNA methylation sites, histone
61 modification markers, gene regulation, and regulatory network [13, 14]. In addition,
62 changes in chromatin accessibility have been implicated with different perspectives of
63 human health as a result of the alterations of nucleosome positioning affected by
64 mutations in chromatin remodelers [15-17].

65 The development of high-throughput sequencing techniques has accumulated a vast
66 amount of chromatin profiles across a variety of cell lines. Large collaborative projects,
67 such as Encyclopedia of DNA Elements (ENCODE) [18], have become a part of the
68 major effort. The Roadmap Epigenomics project provides another similar resource for
69 human stem cells and tissues [19]. Nevertheless, many experimental biologists may
70 lack the bioinformatics expertise to make full use of these valuable resources efficiently.
71 Cistrome DB, a data portal for ChIP-Seq and chromatin accessibility data, although
72 comprises species, factors, biological source, publication, and other information for
73 their collected ChIP-seq and DNase-seq data [20], limited to containing only a part of
74 currently available transcription factors and histone marks. Therefore, it is still very
75 cumbersome and time-consuming to collect, process, and incorporate the chromatin
76 accessibility information of arbitrary genomic regions into bioinformatics and

77 epigenetics studies, which thus makes it difficult to full use of the vast amount of
78 chromatin profiles.

79 We noticed that the epigenome consists of signals from chemical modifications of
80 histones, DNA methylation, non-coding RNA expression, and transcription factors that
81 work in concert to determine the accessibility of the regulatory regions, so-called open
82 region. Then the open regulatory regions can work together with transcription factors,
83 RNA polymerases, and other cellular regulatory machines and produce the final gene
84 expression pattern. In this sense, the chromatin ‘openness’, *i.e.*, the accessibility of
85 genomic regions, bridges the epigenome and transcriptome and plays an important role
86 in understanding the regulatory mechanism. Motivated by the above demand, we built
87 a web server, named OPENANNO, to **annotate** the **openness** of genomic regions across
88 diverse types of cell lines, tissues, and systems. We downloaded the raw sequencing
89 data of 871 DNase-seq experiments across 199 cell lines, 48 tissues and 11 systems
90 from ENCODE data portal [18] and processed by a uniform pipeline. We defined the
91 openness of genomic regions by four statistical strategies, including foreground read
92 count, raw read openness, narrow peak openness, and broad peak openness.
93 Furthermore, we designed a parallel program to enable OPENANNO to efficiently
94 annotate and visualize the openness of a vast amount of genomic regions. We finally
95 demonstrate three applications of OPENANNO in regulatory mechanism and
96 association studies. We believe that this web server will help both the computational
97 and experimental community to facilitate developing methods and discovering
98 important insights, explore the basic biology and various applications, and open a new
99 door to model the regulatory landscape of genome.

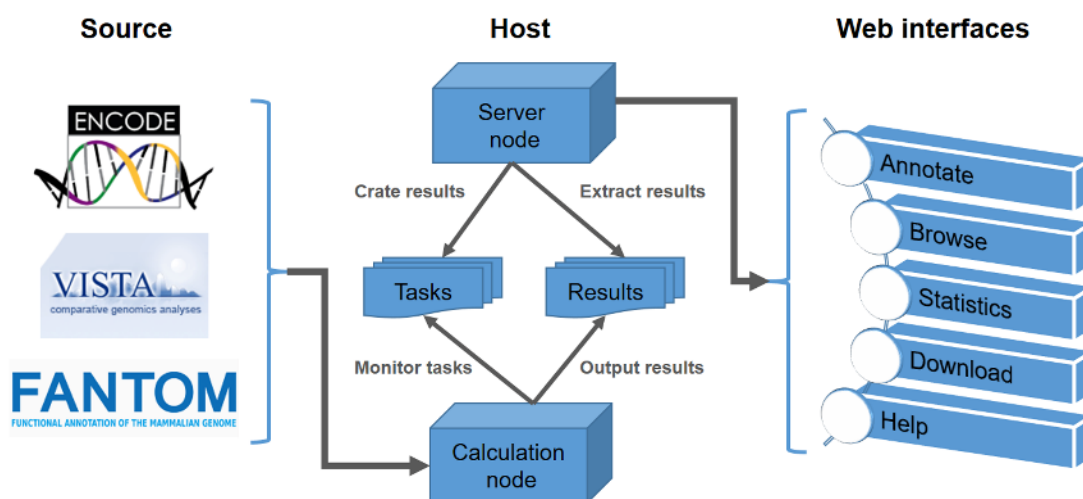
100

101 **Web server content and usage**

102 **Overall design of OPENANNO**

103 As illustrated in **Figure 1**, the diagram for constructing the OPENANNO web server
104 consists of three main parts, *i.e.*, the source, host, and web interface. In the source part,
105 we deposited the meta-information, raw data, and uniformly processed data of 871
106 human DNase-seq experiments from the ENCODE project [18]. It also includes the
107 datasets of regulatory elements that have been experimentally validated from
108 FANTOM[21] and VISTA[22]. The second part, calculation node in the host, monitors
109 the file path that contains annotation tasks, and will automatically calculate and store

110 the results in a specific file path once the server node creates a new task. The server
111 node bridges the host and web interfaces by achieving tasks from users and extracting
112 results or other information for visualization and downloading. The third part, web
113 interface, was developed in a concise and easy-to-use mode. The ‘Annotate’ page
114 provides a service to annotate the openness of a vast amount of genomic regions, and
115 is built on a specially designed infrastructure to extract and visualize big table. The
116 ‘Browse’ page enables users to study the openness of a particular genomic region more
117 intuitively. The ‘Statistics’ page provides detailed information for all the 871 DNase-
118 seq experiments, and an intuitive comparison of the number of experiments in different
119 cell lines, tissues or biological systems. The ‘Download’ page endows users with the
120 ability to directly download the openness of collected experimentally validated
121 regulatory elements. The ‘Help’ page provides other commonly used information to
122 improve the usability of the web server.



123

124 **Figure 1 The diagram for constructing the OPENANNO web server**

125

126 **Web interfaces for annotating**

127 The ‘Annotate’ page, as the major site of OPENANNO, can annotate the openness of
128 genomic regions in batches. As illustrated in **Figure 2**, there are five major steps for
129 this workflow. First of all, we provide a concise task submission approach, which
130 avoids the confusion caused by redundant information to users. By clicking the ‘Browse’
131 button, users can upload a bed or bed.gz file (uncompressed or compressed in gzip
132 format, *e.g.*, ENCF001WKF.bed.gz). The parallel program in our calculation node
133 will extract the first three columns and the sixth column (the chromosomes, starting
134 sites, terminating sites and strands, respectively) separated by tabs for calculating the

135 openness. Note that sorting input in advance using other toolkits, such as bedtools
 136 (<https://bedtools.readthedocs.io>), is preferred for speeding up the calculation. In the
 137 current web server release, we provide the service to annotate the openness of genomic
 138 regions in human reference genome GRCh37 (hg19). We will provide the option of
 139 other genomes or species in future releases. Users can choose to calculate the openness
 140 of these genomic regions in a particular cell line, or directly calculate the openness
 141 all the 871 cell lines. After saving the data as local files, the users can compare the
 142 openness in different cell lines or perform advanced analysis such as the co-openness
 143 analysis, which will be demoed in the following section. Furthermore, users can enable
 144 the option of Per-base pair annotation to calculate the openness of each base-pair of the
 145 genomic regions for some bioinformatics analysis such as machine learning tasks which
 146 will be demoed in the following section.

A Submission form for OPEN ANNO. Fields include: "Select a bed file" (with "Browse" button), "Human GRCh37 (hg19)" (selected), "All cell lines", and "Disable Per Base Pair option". A "Submit" button is at the bottom.

B Progress bar showing 30% completion. A message states: "Your task has been submitted successfully! Uploading file... Initializing the program... Thread loader [00] started. Thread writer [01] started."

E DOWNLOAD dialog box showing a table of files:

FILE	SIZE	MD5
Header file	< 0.01 MB	c5b5e6dd8076f6b461c7b07445b0d8c
Foreground read count	133.43 MB	c06d0189500f06c8526aeb0ca0f73d
Raw read openness	311.24 MB	034390f531d5322a50f682c3721799ec
Narrow peak openness	141.29 MB	85f9a335468292c8b46297369bc762
Broad peak openness	136.13 MB	e7f03e278f79595a029877c20e1e338

D Genomic browser view showing tracks for "RAW read openness, chr11, start 221290, end 221290, strand". It includes tracks for "Chromatin states", "Gene annotations", "Gene expression", "Gene density", "Gene length", "Gene width", "Gene height", "Gene area", "Gene volume", "Gene mass", "Gene weight", "Gene length", "Gene width", "Gene height", "Gene area", "Gene volume", "Gene mass", "Gene weight".

C Results table with columns: "Download", "Raw read openness", "File ID", "Chrom", "Start", "End", "Chr11 98.8%", "Chr12 100.0%", "Chr13 55.2%", and "New task". The table contains 31 rows of genomic data with various numerical values.

148 **Figure 2 Web interfaces for annotating the openness of genomic regions in**
149 **batches**

150 **A.** The interface for submitting a new annotation task. **B.** Display of total calculation
151 progress, and the interface for sending remarks and download links of results by email.
152 **C.** Display of real-time annotation results. Users can browse part of a big table and
153 scroll to any row and any column smoothly. **D.** More detailed information about the
154 experiments and the visualization in UCSC Genome Browser of a specific genomic
155 region. **E.** The interface for downloading the result plain text dump files.

156

157 After submitting the annotation task, users can follow the links of results sent to
158 their email, and download the results after the calculation is completed. Furthermore,
159 users can directly browse real-time results of each region. Here we demonstrate a
160 seamless integration of hardware and software to hold the big result tables which may
161 contain billions of rows and hundreds of columns. Users can observe part of a big table
162 and scroll to any row and any column smoothly. The computational status in different
163 chromosomes will be updated in real time. Users can arbitrarily switch between
164 different chromosomes or the four different types of openness values. By enabling the
165 Color option, each sample in the table will be colored according to the score of openness
166 for intuitively comparing the openness of different genomic regions in different
167 experiments. More detailed information about the experiments and the visualization of
168 a specific genomic region in UCSC Genome Browser [23] can be obtained by double-
169 clicking on a row. After the calculation is completed, users can extract the result files
170 through the ‘Download’ button. By clicking the ‘New Task’ button, users can submit a
171 new task on a new tab in the browser.

172

173 **Web interfaces for browse**

174 On the ‘Browse’ page, users can study the openness of a particular genomic region
175 more intuitively. As illustrated in **Figure 3**, we take an enhancer (chr10: 94,513,996-
176 94,517,989) of the tissue of blood vessels provided by the VISTA Enhancer Browser
177 [22] as an example to demonstrate the service in this page. After submitting the form
178 that contains the chromosome, starting site, terminating site, and strand of a particular
179 genomic region, the web server provides information from three perspectives, including
180 (1) the average openness scores of this genomic region in different biological systems,

181 tissues, and cell lines, (2) the visualization of this genomic region in UCSC Genome
 182 Browser [23], and (3) openness scores and details of the 871 DNase-seq experiments.



183
 184 **Figure 3 The interface for studying the openness of a particular genomic region**
 185 **more intuitively**

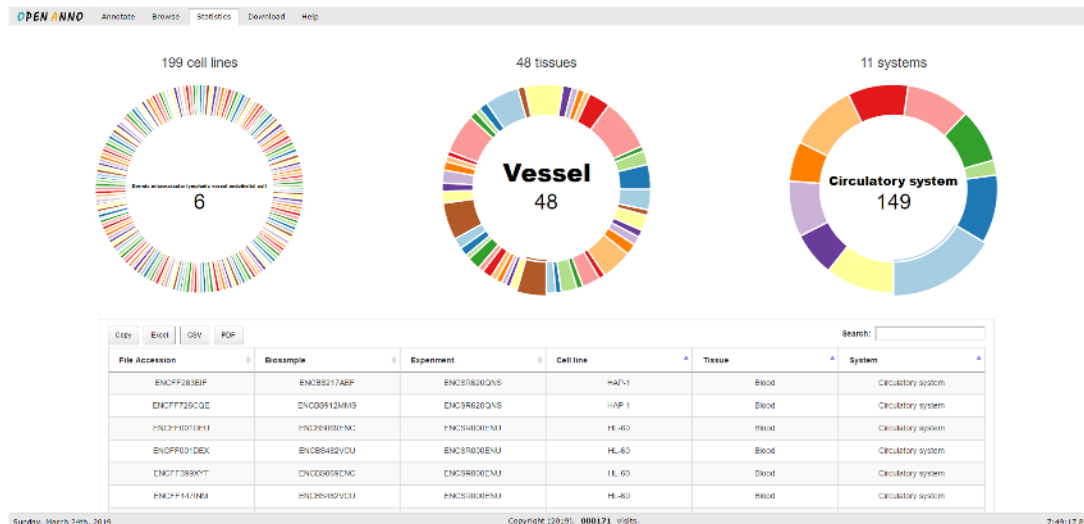
186
 187 First of all, the average openness scores enable users to compare the openness
 188 across different biological systems, tissues, and cell lines more intuitively. For example,
 189 as shown in **Figure 3**, the average openness score of this enhancer in the vessel tissue
 190 is obviously higher than that in other tissues. We calculated z-scores, standard
 191 deviations from the mean, the average openness scores of this enhancer in the vessel
 192 tissue, and found that the z-scores in all the 4 different types of openness definition

193 beyond the three-sigma limits (3.18 of foreground read count, 4.07 of raw read
194 openness, 3.29 of narrow peak openness and 4.37 of broad peak openness), which
195 demonstrates that all the 4 different types of openness of the enhancer in the vessel
196 tissue are significantly higher than that in other tissues from the statistical perspective.
197 This coincides with the fact that this genomic region is an enhancer in blood vessels
198 according to VISTA annotation. Second, the visualization in UCSC Genome Browser
199 provides rapid and reliable displaying of the requested genomic region at any scale,
200 together with dozens of aligned annotation tracks (<http://genome.ucsc.edu/>) [23]. We
201 also provide a hyperlink to UCSC Genome Browser to facilitate users to achieve more
202 concrete information of the particular genomic region. Finally, openness scores and
203 detailed information of the 871 DNase-seq experiments are filled in a table with
204 advanced features. Users can sort the table according to the openness to find out in
205 which experiment the genomic region has higher openness. Users can also sort the table
206 on the basis of different columns according to their own needs. We also provide a
207 convenient service for searching in the table. Users can quickly query information they
208 are interested in. Furthermore, users can directly copy the table or download tables
209 stored in different formats for other requirements.

210

211 **Web interfaces for statistics, downloading, and help**

212 On the ‘Statistics’ page, as shown in **Figure 4**, users can intuitively compare the number
213 of experiments across different cell lines, tissues, or biological systems. A table with
214 advanced features also provides detailed information, including file accessions,
215 biosamples, experiments, cell lines, tissues, and systems, of all the 871 DNase-seq
216 experiments. We collected regulatory element datasets from FANTOM[21] and
217 VISTA[22], including 184,476 FANTOM5 human promoters, 32,693 FANTOM5
218 human enhancers, and 979 VISTA human enhancers. We calculated the openness of
219 these regulatory elements in advance, and provide a download service on the
220 ‘Download’ page to allow users to directly download plain text dump files in
221 compressed (gzip) format. We will continue to provide the openness of other public
222 and validated regulatory elements. To improve the usability, we provide a ‘Help’ page
223 with other commonly used information of the web server, including frequently asked
224 questions, news about the releases of OPENANNO, tutorials of each web interface,
225 citation information, and contact information for help and feedback.



226

227 **Figure 4 The interface for intuitively comparing the number of experiments**
228 **across different cell lines, tissues, or biological systems, and achieving details of all**
229 **the 871 DNase-seq experiments**

230

231 **OPENANNO facilitates regulatory mechanism studies**

232 The chromatin ‘openness’, *i.e.*, the accessibility of genomic regions, calculated using
233 our method has been widely applied to various studies of regulatory mechanism. Here,
234 we show two examples to demonstrate the output of OPENANNO contains valuable
235 information. For example, a model named DeepTACT has been proposed to integrate
236 DNA sequences and chromatin accessibility data for the prediction of chromatin
237 contacts between regulatory elements [24]. Briefly, DeepTACT first performs a one-
238 hot encoding strategy and calculates the raw read openness of each site for
239 characterizing a given genomic region. DeepTACT takes the sequences of two one-hot
240 encoded regulatory elements, and their chromatin openness scores derived from
241 OPENANNO of a given cell type as input. The output of DeepTACT is the predictive
242 score that represents the probability the two regulatory elements have 3D contact. In
243 the deep neural network of DeepTACT, a sequence module is used to extract features
244 from DNA sequences, an openness module is adopted to learn epigenomic features
245 from chromatin openness scores, and an integration module merges outputs of these
246 two modules and extracts high-level features with an attention-based recurrent neural
247 network to predict the probability that the two regulatory elements have 3D contact.

248 Using sequence features and the raw read openness of genomic regions, DeepTACT,
249 as a bootstrapping deep learning model, outperforms existing methods on the task of

250 inferring both promoter-enhancer and promoter-promoter interactions. In more detail,
251 with same test sets, DeepTACT achieves a mean auPRC (the area under the precision-
252 recall curve) score of 0.89 for inferring promoter-promoter interactions compared with
253 0.76 of SPEID [25] and 0.23 of Rambutan [26]. For inferring promoter-enhancer
254 interactions, DeepTACT achieves a mean auPRC of 0.82 compared with 0.67 of SPEID
255 and 0.36 of Rambutan.

256 Besides, DeepTACT provides a finer mapping of promoter-enhancer and promoter-
257 promoter interactions from high-quality promoter capture Hi-C data. Furthermore, the
258 class of hub promoters identified by DeepTACT, and the integrative analysis of existing
259 GWAS data and chromatin contacts predicted by DeepTACT demonstrate the openness
260 calculated by OPENANNO bridges the epigenome and transcriptome and plays an
261 important role in understanding the regulatory mechanism.

262 In addition to DeepTACT, a model named DeepCAPE is proposed to predict
263 enhancers via the integration of DNA sequences and DNase-seq data [27] with the
264 understanding that DNase I hypersensitivity has been shown to be important to identify
265 active cis-regulatory elements including enhancers, promoters, silencers, insulators,
266 and locus control regions [28]. Briefly, DeepCAPE uses the raw read openness of each
267 site of a genomic region as the information of chromatin accessibility to greatly improve
268 the performance of predicting enhancers. In more detail, when the ratio of positive and
269 negative samples is 1:20, the auROC (the area under the receiver operating
270 characteristic curve) and auPRC scores of DeepCAPE are on average 0.151 and 0.590
271 higher than gkmSVM [29], 0.151 and 0.598 higher than DeepSEA [30], and 0.150 and
272 0.588 higher than DeepEnhancer [31], respectively. One-sided paired-sample
273 Wilcoxon signed rank tests consistently suggest that DeepCAPE consistently achieves
274 higher auPRC scores (p -values $< 2.2e-16$ for all the other three baseline methods), and
275 higher auROC scores than a baseline method (p -values $< 2.2e-16$ for all the other three
276 baseline methods).

277 In the model ablation analysis for evaluating the contributions of DNA sequences
278 and DNase-seq data, DeepCAPE illustrates that DNase-seq data provides more
279 information than DNA sequences to greatly improve the performance of prediction.
280 Besides, the information provided by DNA sequences also plays an important role in
281 promoting the performance and making the performance more stable. Because the
282 number of DNase-seq experiments varies between cell lines, the dimensionality of input
283 data varies between cell lines and prevents the use of convolutional neural networks in

284 the cross cell line prediction. DeepCAPE therefore adopts a neural network designed
285 for unsupervised learning of efficient encodings [32], named auto-encoder, to embed
286 chromatin openness scores of a DNA fragment derived from OPENANNO into a vector
287 of fixed length in a low-dimensional latent space. Comparing to the model without an
288 auto-encoder, and other two strategies that average the replicates or randomly select a
289 single replicate, DeepCAPE with an auto-encoder not only makes cross cell line
290 prediction possible, but also maintains superior performance even if the dimensionality
291 of the openness data is reduced. In addition, with a collective scoring strategy,
292 DeepCAPE achieves an average auROC of 0.971 and an average auPRC of 0.862 in
293 the cross cell-line prediction when the ratio of positive and negative samples is 1:20,
294 and thus establishes a landscape of potential enhancers specific to a cell line that still
295 lacks systematic exploration of enhancers.

296 To sum up, DeepCAPE not only achieves superior prediction performance in a cell
297 line-specific manner, but also makes accurate cross cell line predictions possible with
298 the openness scores calculated by OPENANNO. With this understanding, analogous
299 machine learning frameworks can possibly be adapted for the prediction of other
300 functional elements in the genome, including but not limited to promoters, silencers,
301 insulators, repressors, and locus control regions. In addition, the strategy that integrates
302 DNA sequences and chromatin openness can also be generalized for the prioritization
303 of candidate variants in whole-genome sequencing studies, and thus facilitate the
304 regulatory mechanism studies.

305

306 **OPENANNO facilitates association studies**

307 Network-based functional studies play an important role in the identification of disease-
308 associated genes and the interpretation of disease mechanism. The functional
309 relationship of a pair of genes is influenced not only by the co-activation of transcripts,
310 but the regulation mechanism [33]. Besides, the consistence of gene chromatin
311 accessibility indicates the tendency of genes being co-regulated. With this
312 understanding, a gene co-opening network has been constructed based on the raw read
313 openness of genes [34]. Briefly, they take alternative promoters of genes into account
314 when calculating the correlation (absolute value of *Pearson's* correlation coefficient)
315 of the openness scores between two genes, considering the prevalence of the alternative
316 splicing phenomena. By calculating a co-opening score for every pair of genes, they
317 obtain a co-opening matrix for all genes to facilitate the downstream analysis.

318 The results demonstrate that the co-opening network contains new information
319 different from co-expression networks and protein-protein interactions networks. In
320 addition, the genes related to a specific biological process or a specific disease has been
321 demonstrated to tend to be clustered together in the co-opening network, which
322 facilitates detecting functional clusters in the network and predicting new functions for
323 genes. Particularly, through integrative analysis with fruitful genome-wide association
324 studies (GWAS) data, the co-opening network provides a new perspective to the
325 discovery of genes associated with complex diseases, and thus benefits elucidating gene
326 associations and the deciphering of disease mechanisms. For example, by simulating a
327 random walk process on the co-opening network, they use the steady state probability
328 assigned to a gene as a score to measure the likelihood that the gene is associated with
329 the disease under investigation. Applying this strategy to a complex disease named
330 *Psoriasis*, a potentially disfiguring immune-mediated inflammatory disease of skin,
331 they discovered *TNFSF14* (TNF super family member 14, a biomarker of *Psoriasis*
332 [35]), which was ranked second by the random walk model while cannot be identified
333 by GWAS (p -value = 0.1616, ranked 1259 based on the p -value). In general, the co-
334 opening network is ready to serve as a useful resource complementary to the widely
335 used co-expression network, and thus shed light on the studies in system biology.

336

337 **Perspectives and concluding remarks**

338 Chromatin accessibility, which bridges the epigenome and transcriptome, is a very
339 valuable resource for interpreting non-coding genomic region and understanding the
340 regulatory mechanism. In this study, we downloaded raw sequencing data of 871
341 DNase-seq experiments across 199 cell lines, 48 tissues and 11 biological systems from
342 ENCODE, and defined the openness of genomic regions from four perspectives. In
343 addition, we take an enhancer of the tissue of blood vessels provided by the VISTA
344 Enhancer Browser as an example to illustrate the valuable information provided by all
345 the four different types of openness from statistical perspective. Furthermore, we
346 designed a parallel program to endow OPENANNO with the ability to efficiently
347 annotate and visualize openness for a vast amount of genomic regions. Finally, we
348 introduced three examples to demonstrate the output of OPENANNO serves as
349 valuable input for follow-up regulatory mechanism and association studies.

350 Our web server has four main application scenarios. First, one can use our web
351 server to annotate openness of genomic regions, and then integrate the information of
352 openness to a machine learning model for superior performance. Second, one can use
353 our web server to visualize the openness of a specific genomic region to intuitively
354 understand this region has higher openness in which cell lines, tissues, or systems, and
355 thus contribute to the study of functional implications of this genomic region. Third,
356 our web server offers a new opportunity to reinterpret abundant data cumulated by
357 genome-wide association studies, and thus one can characterize variants by integrating
358 upstream openness annotated with our web server and downstream gene expression.
359 Finally, one can use the openness annotated with our web server to construct gene co-
360 opening networks which provide a new perspective to association studies.

361 To better serve the academic community, we will continue to collect public data
362 and update OPENANNO regularly in the future. Our next plan is to provide the option
363 of other genomes or species, and the option of annotating using other chromatin
364 accessibility data, such as ATAC-seq data. We will continue to provide the openness
365 of other public and validated regulatory elements for downloading directly. According
366 to users' feedbacks, we will continue to improve the interfaces and performance of
367 OPENANNO. We believe that OPENANNO would serve as a useful tool for both
368 bench scientists and computational biologists, and shed light on studies including but
369 not limited to bioinformatics and system biology.

370

371 **Materials and methods**

372 **Data collection**

373 We first parsed a total of 41,418 JSON files from ENCODE to obtain detailed
374 information about experiments, biosamples, cell lines, tissues, and systems of the
375 DNase-seq data provided by the ENCODE project [18]. Under the constraint that each
376 experiment contains both narrow peaks and broad peaks, we downloaded raw
377 sequencing data of 891 DNase-seq experiments in human reference genome GRCh37
378 (hg19), and then identified experiments corresponding to 199 cell lines, 48 tissues and
379 11 biological systems. We collected datasets of regulatory elements from FANTOM
380 [21] and VISTA[22] that have been experimentally validated, including 184,476
381 FANTOM5 human promoters, 32,693 FANTOM5 human enhancers, and 979 VISTA
382 human enhancers.

383

384 **Definition of openness**

385 We defined the openness of given genomic regions from four perspectives, including
386 foreground read count, raw read openness, narrow peak openness, and broad peak
387 openness. Specifically, given the raw sequencing data of a DNase-seq experiment, we
388 provided the number of reads (N), *i.e.*, foreground read count, falling at a specific
389 genomic region to facilitate special applications that may use raw read counts of a
390 DNase-seq experiment directly. To remove the effect of sequencing depth, we defined
391 the raw read openness (S) of a genomic region as the foreground read count (N) divided
392 by the average number of reads falling at a position in a background region of size W
393 surrounding the given region. The raw read openness (S) can be simply calculated as

$$S = \frac{N}{K/W} \quad (1),$$

394 where K is the number of reads falling into the background region of size W . The size
395 of a background region W is set to 1 M base pairs, according to the suggestion from
396 [36].

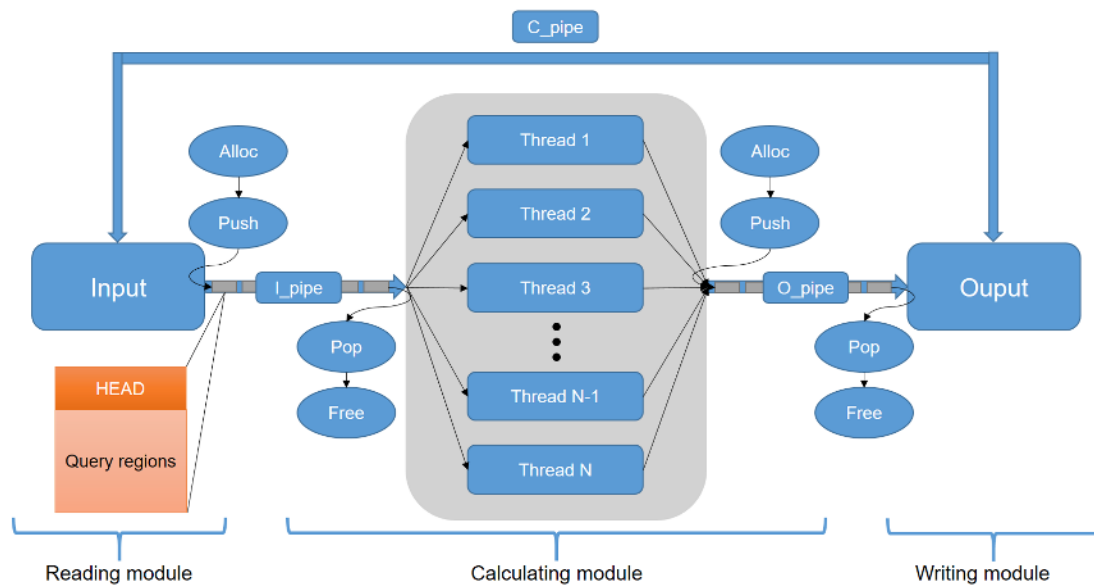
397 Analogously, we defined the narrow peak openness (and broad peak openness) of a
398 genomic region as the number of narrow peaks (or broad peaks) overlapping with the
399 genomic region, divided by the average number of narrow peaks (or broad peaks)
400 overlapping with a position in a background region of size W surrounding the given
401 genomic region.

402

403 **Parallel computing and real-time browsing**

404 To facilitate querying reads and peaks at high frequencies, and the massive demand for
405 annotating openness of a large number of genomic regions, we designed a parallel
406 strategy that endows OPENANNO with an ability to efficiently annotate openness of a
407 vast amount of genomic regions. We used C++, a programming language known for its
408 high efficiency, to develop a multithreaded program that consists of a reading module,
409 a calculating module, and a writing module. As illustrated in **Figure 5**, the reading
410 module packages the query regions in the input file into data blocks and pushes them
411 into the input pipeline I_pipe . Each data block contains a portion of the query regions,
412 and header information that contains the number and indexes of these query regions in
413 the input file. When a thread in the calculating module is idle, the data block is
414 automatically extracted from the input pipeline for calculation. The results are pushed

415 into the output pipeline O_pipe in a similar form of data blocks, which are then popped
416 out and written into the disks by the writing module. Through the communication
417 pipeline C_pipe , the reading module can respond to the working state of the writing
418 module. When the writing speed is lower than the reading speed, the reading module
419 pauses pushing the data block into I_pipe to effectively save memory resources. With
420 the parallel program, OPENANNO can calculate the four types of openness across all
421 the 871 DNase-seq experiments of one thousand genomic regions within 1 second.



422

423 **Figure 5 The multithreaded program for efficiently annotating openness of a**
424 **vast amount of genomic regions**

425

426 Users may need to annotate a large number of genomic regions, and thus results in
427 very large tables of openness, especially when the per-base option is enabled. Loading
428 all the results directly to the front-end of OPENANNO will take a long time and result
429 in an unsmooth experience, and thus is an unrealistic choice. It is, therefore, necessary
430 to provide a highly efficient front-end web application to browse these tables. To
431 accomplish this, we used WebSocket (<https://www.websocket.org/>), a computer
432 communications protocol, to provide full-duplex communication channels over a single
433 socket over the web and the remote host, realize the real-time browsing of calculation
434 results, and thus users can observe part of a large table and scroll to any row and any
435 column smoothly.

436

437 **Web server implementation**

438 The whole design of the OPENANNO is shown in Figure 1, with the possible jumps
439 among web pages illustrated. OPENANNO is freely available to all users without a
440 login requirement. The current version of OPENANNO was deployed on a calculation
441 node and a server node of a high-performance computer cluster. The calculation node,
442 whose operating system is CentOS 7.5 (one of the most popular Linux distributions;
443 <https://www.centos.org>), has 56 hyper-threaded processors to perform efficient parallel
444 computing, RAM of 775 GB to support a large number of reading and writing
445 operations, and storage space of 321 TB to store the vast amounts of data. Incron
446 (<http://inotify.aiken.cz/?section=incron>) is used to monitor filesystem events and
447 executes predefined commands. Once the server node creates a new task in a specific
448 file path, the parallel program automatically runs and stores the results in another
449 specific file path for the website to read and users to download.

450 In the server node, a Linux-based Nginx server (<https://www.nginx.com/>), which
451 uses dramatically less memory than Apache (<https://www.apache.org/>) and can handle
452 roughly four times more requests per second, is deployed to improve the performance,
453 reliability, and security of OPENANNO. PHP (<http://www.php.net/>) is used for server-
454 side scripting. Bootstrap (a popular toolkit for developing with HTML, CSS and
455 JavaScript; <https://getbootstrap.com/>) and jQuery (a fast and feature-rich library
456 designed to simplify the JavaScript programming; <https://jquery.com/>) are adopted for
457 the front-end, *i.e.*, the interactive and responsive user interface, of OPENANNO. The
458 user interface of OPENANNO can automatically respond to the devices and browsers
459 with different screen resolutions, and change its structure and shape according to the
460 resolution in order to optimize the visualization. DataTables (<https://datatables.net/>), as
461 a plug-in for the jQuery and JavaScript library, is used to add advanced features to the
462 tables. The visualizations of bar charts and pie charts are implemented with JavaScript
463 libraries named CanvasJS (<https://canvasjs.com/>) and morris.js
464 (<https://morrisjs.github.io/morris.js/index.html>), respectively. To obtain the
465 visualization in UCSC Genome Browser of a query region, we put the information of
466 chromosome, start position, end position and reference genome to the link
467 [http://genome.ucsc.edu/cgi-](http://genome.ucsc.edu/cgi-bin/hgTracks?db=(reference_genome)&position=(chromosome):(start_position)-(end_position))
468 [bin/hgTracks?db=\(reference_genome\)&position=\(chromosome\):\(start_position\)-](http://genome.ucsc.edu/cgi-bin/hgTracks?db=(reference_genome)&position=(chromosome):(start_position)-(end_position))
469 [\(end_position\)](http://genome.ucsc.edu/cgi-bin/hgTracks?db=(reference_genome)&position=(chromosome):(start_position)-(end_position)) to finish UCSC link construction.

470 All codes are developed using Vim (a highly flexible text editor that supports any
471 kind of text; <https://www.vim.org/>). The performance of OPENANNO has been tested
472 in Chrome, Firefox, Opera and Microsoft Edge on Windows 10, Ubuntu 16.04 and
473 MacOS 10.12. We hope users could feedback their comments and suggestions through
474 the contact page on our website to help us improve OPENANNO.
475

476 **Authors' contributions**

477 RJ and YW designed the project. SC and RJ collected data and implemented the web
478 server. SC, YW and RJ wrote the paper. All authors read and approved the final
479 manuscript.

480

481 **Competing interests**

482 The authors have declared no competing interests.

483

484 **Acknowledgments**

485 This work was partially supported by the National Key Research and Development
486 Program of China (Grant No. 2018YFC0910404), the National Natural Science
487 Foundation of China (Grant Nos. 61873141, 61721003, 61573207, 11871463, and
488 61671444), and the Tsinghua-Fuzhou Institute for Data Technology. Rui Jiang is a
489 RONG professor at the Institute for Data Science, Tsinghua University.

490

491 **Reference**

- 492 [1] Frazer KA. Decoding the human genome. *Genome Res* 2012;22(9):1599-601.
493 [2] Sosnay PR, Cutting GR. Interpretation of genetic variants. *Thorax* 2014;69(3):295-
494 7.
495 [3] Klemm SL, Shipony Z, Greenleaf WJ. Chromatin accessibility and the regulatory
496 epigenome. *Nat Rev Genet* 2019;20(4):207-20.
497 [4] Tsompana M, Buck MJ. Chromatin accessibility: a window into the genome.
498 *Epigenetics Chromatin* 2014;7(1):33.
499 [5] Radman-Livaja M, Rando OJ. Nucleosome positioning: how is it established, and
500 why does it matter? *Dev Biol* 2010;339(2):258-66.
501 [6] Rizzo JM, Sinha S. Analyzing the global chromatin structure of keratinocytes by
502 MNase-seq. *Methods Mol Biol* 2014;1195:49-59.
503 [7] Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. Transposition of
504 native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-
505 binding proteins and nucleosome position. *Nat Methods* 2013;10(12):1213-8.
506 [8] Cui K, Zhao K. Genome-wide approaches to determining nucleosome occupancy
507 in metazoans using MNase-Seq. *Methods Mol Biol* 2012;833:413-9.
508 [9] Giresi PG, Kim J, McDaniell RM, Iyer VR, Lieb JD. FAIRE (Formaldehyde-
509 Assisted Isolation of Regulatory Elements) isolates active regulatory elements from
510 human chromatin. *Genome Res* 2007;17(6):877-85.

- 511 [10]John S, Sabo PJ, Canfield TK, Lee K, Vong S, Weaver M, et al. Genome-scale
512 mapping of DNase I hypersensitivity. *Curr Protoc Mol Biol* 2013;Chapter 27:Unit 21
513 7.
- 514 [11]Simon JM, Giresi PG, Davis IJ, Lieb JD. Using formaldehyde-assisted isolation of
515 regulatory elements (FAIRE) to isolate active regulatory DNA. *Nat Protoc*
516 2012;7(2):256-67.
- 517 [12]John S, Sabo PJ, Thurman RE, Sung MH, Biddie SC, Johnson TA, et al. Chromatin
518 accessibility pre-determines glucocorticoid receptor binding patterns. *Nat Genet*
519 2011;43(3):264-8.
- 520 [13]Ward LD, Kellis M. Evidence of abundant purifying selection in humans for
521 recently acquired regulatory functions. *Science* 2012;337(6102):1675-8.
- 522 [14]Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, Shendure J. A general
523 framework for estimating the relative pathogenicity of human genetic variants. *Nat*
524 *Genet* 2014;46(3):310-5.
- 525 [15]Gaspar-Maia A, Alajem A, Polesso F, Sridharan R, Mason MJ, Heidersbach A, et
526 al. Chd1 regulates open chromatin and pluripotency of embryonic stem cells. *Nature*
527 2009;460(7257):863-8.
- 528 [16]Hargreaves DC, Crabtree GR. ATP-dependent chromatin remodeling: genetics,
529 genomics and mechanisms. *Cell Res* 2011;21(3):396-420.
- 530 [17]Schwartzentruber J, Korshunov A, Liu XY, Jones DT, Pfaff E, Jacob K, et al.
531 Driver mutations in histone H3.3 and chromatin remodelling genes in paediatric
532 glioblastoma. *Nature* 2012;482(7384):226-31.
- 533 [18]Consortium EP. An integrated encyclopedia of DNA elements in the human
534 genome. *Nature* 2012;489(7414):57-74.
- 535 [19]Bernstein BE, Stamatoyannopoulos JA, Costello JF, Ren B, Milosavljevic A,
536 Meissner A, et al. The NIH Roadmap Epigenomics Mapping Consortium. *Nat*
537 *Biotechnol* 2010;28(10):1045-8.
- 538 [20]Mei S, Qin Q, Wu Q, Sun H, Zheng R, Zang C, et al. Cistrome Data Browser: a
539 data portal for CHIP-Seq and chromatin accessibility data in human and mouse. *Nucleic*
540 *Acids Res* 2017;45(D1):D658-D62.
- 541 [21]Lizio M, Harshbarger J, Shimoji H, Severin J, Kasukawa T, Sahin S, et al.
542 Gateways to the FANTOM5 promoter level mammalian expression atlas. *Genome Biol*
543 2015;16:22.
- 544 [22]Visel A, Minovitsky S, Dubchak I, Pennacchio LA. VISTA Enhancer Browser--a
545 database of tissue-specific human enhancers. *Nucleic Acids Res* 2007;35(Database
546 issue):D88-92.
- 547 [23]Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, et al. The
548 human genome browser at UCSC. *Genome Res* 2002;12(6):996-1006.
- 549 [24]Li W, Wong WH, Jiang R. DeepTACT: predicting 3D chromatin contacts via
550 bootstrapping deep learning. *Nucleic Acids Res* 2019.
- 551 [25]Singh S, Yang Y, Póczos B, Ma J. Predicting Enhancer-Promoter Interaction from
552 Genomic Sequence with Deep Neural Networks. *bioRxiv* 2018:085241.

- 553 [26]Schreiber J, Libbrecht M, Bilmes J, Noble W. Nucleotide sequence and DNaseI
554 sensitivity are predictive of 3D chromatin architecture. *bioRxiv* 2018:103614.
- 555 [27]Chen S, Gan M, Lv H, Jiang R. DeepCAPE: a deep convolutional neural network
556 for the accurate prediction of enhancers. *bioRxiv* 2018:398115.
- 557 [28]Thurman RE, Rynes E, Humbert R, Vierstra J, Maurano MT, Haugen E, et al. The
558 accessible chromatin landscape of the human genome. *Nature* 2012;489(7414):75-82.
- 559 [29]Lee D. LS-GKM: a new gkm-SVM for large-scale datasets. *Bioinformatics*
560 2016;32(14):2196-8.
- 561 [30]Zhou J, Troyanskaya OG. Predicting effects of noncoding variants with deep
562 learning-based sequence model. *Nat Methods* 2015;12(10):931-4.
- 563 [31]Min X, Zeng W, Chen S, Chen N, Chen T, Jiang R. Predicting enhancers with deep
564 convolutional neural networks. *BMC Bioinformatics* 2017;18(Suppl 13):478.
- 565 [32]Hinton GE, Salakhutdinov RR. Reducing the dimensionality of data with neural
566 networks. *Science* 2006;313(5786):504-7.
- 567 [33]Hecker M, Lambeck S, Toepfer S, van Someren E, Guthke R. Gene regulatory
568 network inference: data integration in dynamic models-a review. *Biosystems*
569 2009;96(1):86-103.
- 570 [34]Li W, Wang M, Sun J, Wang Y, Jiang R. Gene co-opening network deciphers gene
571 functional relationships. *Mol Biosyst* 2017;13(11):2428-39.
- 572 [35]Chandran V, Cook RJ, Edwin J, Shen H, Pellett FJ, Shanmugarajah S, et al. Soluble
573 biomarkers differentiate patients with psoriatic arthritis from those with psoriasis
574 without arthritis. *Rheumatology (Oxford)* 2010;49(7):1399-405.
- 575 [36]Duren Z, Chen X, Jiang R, Wang Y, Wong WH. Modeling gene regulation from
576 paired expression and chromatin accessibility data. *Proc Natl Acad Sci U S A*
577 2017;114(25):E4914-E23.
- 578

579 **Figure legends**

580 **Figure 1 The diagram for constructing the OPENANNO web server**

581

582 **Figure 2 Web interfaces for annotating the openness of genomic regions in**
583 **batches**

584 **A.** The interface for submitting a new annotation task. **B.** Display of total calculation
585 progress, and the interface for sending remarks and download links of results to email.

586 **C.** Display of real-time annotation results. Users can observe part of a big table and
587 scroll to any row and any column smoothly. **D.** More detailed information about the
588 experiments and the visualization in UCSC Genome Browser of a specific genomic
589 region. **E.** The interface for downloading the result plain text dump files.

590

591 **Figure 3 The interface for studying the openness of a particular genomic region**
592 **more intuitively**

593

594 **Figure 4 The interface for intuitively comparing the number of experiments**
595 **across different cell lines, tissues or systems, and achieving detail information of**
596 **all the 871 DNase-seq experiments**

597

598 **Figure 5 The multithreaded program for efficiently annotating openness of a**
599 **vast amount of Genomic regions**