# 1 ChIP-AP – An Integrated ChIP-Seq Analysis Pipeline

2

3 **Jeremiah Suryatenggara[1], Kol Jia Yong[1,2], Danielle E. Tenen[3], Daniel G. Tenen[1,4], Mahmoud A.**
4 **Bassal[1,4]**

5 1 – Cancer Science Institute of Singapore, National University of Singapore, Singapore

6 2 – Department of Biochemistry, Yong Loo Lin School of Medicine, National University of
7 Singapore, Singapore

8 3 – Broad Institute of MIT and Harvard, Boston, Massachusetts, USA

9 4 – Harvard Stem Cell Institute, Boston, Massachusetts, USA

10

11 **Correspondence:**
12 daniel.tenen@nus.edu.sg
13 mahmoud.bassal@mymail.unisa.edu.au

14

## 15 Abstract

16 ChIP-Seq is a technique used to analyse protein-DNA interactions. The protein-DNA complex is

17 pulled down using a protein antibody, after which sequencing and analysis of the bound DNA

18 fragments is performed. A key bioinformatics analysis step is "peak" calling - identifying regions of

19 enrichment. Benchmarking studies have consistently shown that no optimal peak caller exists. Peak

20 callers have distinct selectivity and specificity characteristics which are often not additive and seldom

21 completely overlap in many scenarios. In the absence of a universal peak caller, we rationalized one

22 ought to utilize multiple peak-callers to 1) gauge peak confidence as determined through detection by

23 multiple algorithms, and 2) more thoroughly survey the protein-bound landscape by capturing peaks

24 not detected by individual peak callers owing to algorithmic limitations and biases. We therefore

25 developed an integrated ChIP-Seq Analysis Pipeline (ChIP-AP) which performs all analysis steps

26 from raw fastq files to final result, and utilizes four commonly used peak callers to more thoroughly

27 and comprehensively analyse datasets. Results are integrated and presented in a single file enabling

28 users to apply selectivity and sensitivity thresholds to select the consensus peak set, the union peak

29 set, or any sub-set in-between to more confidently and comprehensively explore the protein-bound

30 landscape. (https://github.com/JSuryatenggara/ChIP-AP).

## Introduction

ChIP-Seq is an extensively used experimental technique that aims to identify DNA binding location sequences and motifs of DNA-interacting proteins such as transcription factors[1,2], histone-modifier proteins[3], or novel DNA-binding proteins. To perform a ChIP-Seq experiment, cells are fixed, the chromatin-protein complex sonicated, and the DNA fragments interacting with the protein of interest pulled down by the targeted protein antibody. Following experimental pull-down of the protein-bound DNA and sequencing, the raw sequencing data (raw fastq files) undergoes processing and analysis, after which, biological relevance can be inferred[4]. Such analyses, in conjunction with follow up mechanistic studies, shed light on the DNA-associated proteins biological function and roles[4].

Since the development of ChIP-Seq[2,5,6], computational analysis of ChIP-Seq experiments has always been a multi-step process requiring multiple command line programs which, to use most effectively, requires knowledge and experience in computing and programming[4]. Wet-lab biologists without command line or coding experience have typically relied on bioinformaticians, with their computing expertise, to analyse the sequenced data despite them potentially having a reduced understanding of the underlying biology. Irrespective of whom performs the analysis however, in the computational space, many analogous programs have been developed for each stage of the analysis, complicating how one should approach an analysis and the decision of which programs to use in conjunction with one-another. It is therefore easy to understand why two different analysis methodologies, even if they appear superficially identical or similar, will almost certainly report different results, leading to conflicting conclusions from the same biological experiment[7,8]. To complicate matters further, published methods outlining the workflow used to analyse a dataset will consistently lack essential details, with some authors omitting key program modification parameters/flags, or neglecting to include key analysis steps entirely, relegating published analyses to being almost entirely irreproducible for other researchers.

58      Of all the complications plaguing ChIP-Seq analyses though, perhaps the most well-known source of

59      inconsistency between analyses is the choice of peak calling algorithm[2,7,8]. This observation was

60      convincingly demonstrated by Steinhauser et. al.[8] in their study comparing 20 peak callers wherein

61      they reported poor agreement between the called peak sets across the profiled callers. This, and other

62      studies, therefore show that peak callers have distinct selectivity and specificity characteristics which

63      are often not additive and seldom completely overlap in many scenarios[8-11]. Consequently, such

64      differing operating characteristics results in a lack of consistency across the reported regions of

65      enrichment (and associated genes) for each peak caller. This has follow-on effects in that

66      downstream functional analysis for the protein of interest would therefore give differing, potentially

67      conflicting results. Additionally, it has been shown extensively that the performance of a peak caller

68      is subject to the read characteristics and read distributions of the dataset in question[8-12]. An individual

69      peak caller can outperform other callers in certain datasets but will perform poorly in alternate

70      datasets. Therefore, relying on a one-caller-fits-all approach when analysing datasets with different

71      DNA-binding proteins, immunoprecipitation and library preparation protocols is objectively not the

72      soundest approach to yield reliable, consistent and comprehensive results.

73

74      We therefore rationalized that in order to improve the reliability and consistency of a ChIP-Seq

75      analysis, one ought to focus on and improve the consistency, confidence and comprehensiveness of

76      the peak detection step without requiring additional wet-lab observations. To address this, we

77      designed our ChIP-Seq Analysis Pipeline (ChIP-AP) which integrates all processes of a ChIP-Seq

78      analysis (from raw fastq to final result) into a single, easy to use package, that utilizes four

79      commonly used peak callers[13-17] (for either transcription factors or histone modifier proteins), and

80      integrates their results into a single output file, from which, users are able to infer peak confidence. If

81      a peak is called by multiple callers, one can infer that the reported peak has a higher confidence and

82      is less likely to be a false-positive or an artefact of algorithmic bias or limitation. Alternatively, by

83      integrating the results of all the callers (the union of all peaks), one can more comprehensively

84      survey the binding landscape of the binding protein by capturing peaks that would otherwise be

85      uncalled or "lost" if relying on a lone peak caller. In other words, the union peak set enables a more

86     comprehensive survey of the binding landscape by accepting all peaks irrespective of confidence. By

87     utilizing multiple peak callers and integrating their results, users are able to determine the selectivity

88     and specificity requirements that best describe their dataset while allowing them to circumvent

89     inherent sample characteristics that result in poor peak calling performance of any single peak caller,

90     thus enabling the capture of either the most number or peaks (the union of all peaks), the most

91     confident peaks (the consensus results), or any sub-set of the gradient in-between. ChIP-AP can

92     therefore become an effective tool for users by providing both substantial improvements to peak

93     capturing and analysis reliability. ChIP-AP is available on GitHub

94     (https://github.com/JSuryatenggara/ChIP-AP) with extensively detailed wiki pages describing

95     installation, use and results interpretation (https://github.com/JSuryatenggara/ChIP-AP/wiki).

96

## Results

### Consensus peaks increase motif and ontology accuracy

99     A reproducible result instils greater confidence in its validity. Likewise, ChIP-Seq peaks that can be

100     detected by multiple peak callers, each utilizing different peak detection algorithms, garner greater

101     confidence than peaks called by an individual peak caller. ChIP-AP, with its utilization of four

102     different peak callers, reports along-side the coordinates of a peak how many callers detected the said

103     peak. Using this, users are able to filter on the consensus peaks, which are the peaks detected by all

104     four peak callers and, consequently, carry the greatest confidence. We hypothesized that utilizing the

105     consensus peak set would increase the percentage of peaks containing a valid binding motif

106     (peak-motif percentage) without drastically affecting the motif position bias (the distance of the motif

107     to the weighted peak center). The consensus peak set would also increase the likelihood of identifying

108     correct binding motifs while masking co-factor binding artifacts. Finally, the consensus peak set

109     should also improve gene ontology (GO) results by ensuring only the strongest binding candidates are

110     included in analysis.

111

112     To investigate whether the peak-motif percentage for the consensus peak set is significantly better

113     than using a single peak caller (MACS2), we processed 10 transcription factor (TF) datasets from

114    differing TF families, across 3 cell lines sourced from ENCODE[18], and determined their peak-motif

115    percentage (see **Methods** for sample details). For each TF, we downloaded the binding motif from

116    MethMotif[19] (which sourced its motifs from JASPAR 2018[20]), and determined how many peaks

117    contain the binding motif while allowing up to a single sequence mismatch. As expected, a significant

118    difference in the average peak-motif percentage was observed across all 10 TF's (two-tailed t-test

119    p=0.0012, **Figure 1a**). The degree by which the consensus peak set improved the peak-motif

120    percentage was variable, with an up to 90% improvement for RUNX1, but none the less, still showed

121    improvement for all 10 TFs over the MACS2 peak set alone (**Supplemental Table 1**).

122

123    We next investigated whether the consensus peak set significantly altered the motif position bias with

124    respect to the weighted peak center as compared to individual callers. Across all 10 TFs, the benefit of

125    the consensus peak set was variable, with it out-performing all individual peak callers in some

126    datasets (**Figure 1b, Supplemental Table 2**), while in others, providing comparable results

127    (**Supplemental Figure 1a, Supplemental Table 2**). Consistently though, the consensus peak set did

128    show significant improvement over at least half the peak callers tested suggesting that utilizing the

129    consensus peak set will either give an improved motif position bias profile or report comparable

130    results to having used an individual peak caller.

131

132    Next, we questioned whether the consensus peak set can provide improved *de novo* motif sequence

133    detection while masking co-factor binding artefacts. Previously, Lin et al.[21] described that *de novo*

134    global motif analyses can potentially be contaminated by co-factor motif sequence artefacts. In their

135    publication, they highlighted potential co-factor motif artefacts for CEBPB and MAFF. Using

136    ChIP-AP, we performed *de novo* motif analysis using HOMER[15] and the MEME-Suite[22-24] for the

137    MACS2 called and the consensus peak set for these two TFs. For CEBPB, the first candidate motif hit

138    reported by HOMER for both the consensus and the MACS2 peak sets was near identical, which,

139    according to Lin et al.[21], contains co-factor binding motif artefacts (**Figure 1c, d** upper panels).

140    However, the second motif hit for the consensus peak set (**Figure 1c** lower panel) shows the clean

141    binding motif sequence TTGC, which is the CEBPB motif that contains neither co-factor motif

142    contamination nor a heterodimer sequence[21]. The MACS2 peak set's second motif result however

143    (**Figure 1d**, lower panel), failed to report the same motif result. When analysed using the

144    MEME-Suite, the consensus peak set showed two motifs with co-factor artefacts (**Supplemental**

145    **Figure 1b**, 1st and 3rd ranked), two motifs with the heterodimer sequence (**Supplemental Figure 1b**,

146    2nd and 4th ranked) and the fifth result showed the TTGC binding motif without co-factor motif

147    contamination nor a heterodimer sequence. Conversely, for the MACS2 dataset as analysed by the

148    MEME-Suite, three motif results contain co-factor sequence artefacts (**Supplemental Figure 1c**, 1st,

149    3rd and 4th ranked), one result with heterodimer sequence (**Supplemental Figure 1c**, 2nd ranked) and

150    the fifth result was also the TTGC binding motif without co-factor motif contamination nor a

151    heterodimer sequence. Therefore, for CEBPB, although both *de novo* motif algorithms reported

152    similar findings, the consensus peak set showed a cleaner and more direct signal for the CEBPB

153    binding motif from the second HOMER result, a finding not immediately evident in the MACS2 set

154    without careful inspection of the data.

155

156    For the TF MAFF, the HOMER *de novo* motif result for the consensus peak set remained more

157    consistent (**Figure 1e**) with both the top two motif results showing the binding sequence TCAGCA.

158    The MACS2 peak set however, wasn't as consistent showing different sequences between the first

159    and second motif candidates (**Figure 1f**). Similarly for the MEME-Suite results (**Supplemental**

160    **Figure 1d, e**), MEME consistently calls the TGCTGA for both peak sets but the 6th reported motif

161    candidate for the consensus peak set shows a heterodimer binding profile (characterised by the

162    "sequence - spacer - reverse complement sequence" profile), a result not recapitulated in the MACS2

163    peak set entirely, thereby supporting the notion that the consensus peak set can provide more direct

164    *de novo* enriched motif results over using a single peak caller alone.

165

166    Our final investigation was to test whether the consensus peak set can provide improved, more direct

167    gene ontology (GO) results. In running a GO analysis for all 10 TF's profiled and comparing the

168    consensus peak set results to the MACS2 peak sets, we observed that for certain datasets, such as

169    RUNX1, ATF4, JUN, ZBTB33 and GATA1, the consensus peak set GO results returned more

170    relevant and directly related terms than the MACS2 peak set (**Table 1, Supplemental Tables 3, 4**).

171    For the RUNX1 results, whereas the top 20 MACS2 GO results contained generic GO terms, the

172    consensus peak GO listing clearly outlined RUNX1 functions regarding hematopoietic differentiation,

173    regulation of metabolic and signalling pathways and autophagy regulation, all of which are known

174    published functions of RUNX1[25-30] (**Table 1**). The GO terms returned when searching the consensus

175    peak candidate list can therefore, for certain datasets, provide significantly clearer and more direct GO

176    results by providing information on only the gene terms corresponding to the most confident peaks

177    called by all peak callers.

178

179    Therefore, utilizing the consensus peak set can provide added benefits to identify novel binding motifs

180    or to identify more direct biological processes modulated by a protein of interest, especially if it is not

181    well characterized as evidenced by the results presented. In all metrics investigated, the consensus

182    peak set's performance was either significantly improved, or, in worst performing cases, provided

183    results comparable to having used only a single peak caller.

184

185    Capturing Lost Peaks with the Union Peak Set
186    A number of variables can affect a ChIP-Seq experiments efficiency resulting in poor enrichment and

187    potentially giving rise to a high signal:noise ratio dataset. Every peak caller has differing operating

188    characteristics and thus, has differing abilities to handle these difficult to process datasets[8,12]. A

189    ChIP-Seq analysis utilizing only a single peak caller would be solely dependent on the chosen peak

190    callers' ability to handle the signal:noise ratio and enrichment characteristics of that dataset. If the

191    peak caller struggles to differentiate signal from noise effectively, few peaks will be called and a

192    dataset will give an inconclusive result owing to its ineffectiveness to deal with the dataset. However,

193     some peak callers are more capable at handling difficult datasets, and so an experiment may show

194     poor enrichment, but it simply needs to be analysed with the right peak caller for its specific

195     characteristics, the choice of which may not be evident or obvious in advance.

196

197     One protein that is relatively difficult to perform ChIP-Seq on, is the oncogene sal-like protein 4

198     (SALL4). SALL4 has been shown to play essential roles in maintaining pluripotency and self-renewal

199     characteristics of embryonic stem cells (ESC)[31]. It is typically down-regulated after birth but has been

200     found to be aberrantly regulated in many tumors[31,32]. Studies have also shown SALL4 to have

201     multiple protein interacting partners and DNA-binding and regulation functions[31,33]. An attempt to

202     capture the DNA-binding partners of SALL4 was undertaken with the sequenced result showing poor

203     enrichment on the fingerprint plot with little separation between the SALL4 ChIP-Seq replicate and

204     control curves (**Figure 2a**), indicating it will likely be difficult to call peaks for this dataset. When

205     processed with ChIP-AP, we observed that peak callers GEM and MACS2 struggle to call peaks

206     (**Figure 2b**) with each returning a total of 1,362 and 1,937 peaks respectively. HOMER, is able to call

207     approximately double the number of peaks at 3,760. However, Genrich, which determines peaks using

208     an area under the curve (AUC) calculation rather than generating a Poisson distribution model (as

209     seen in MACS2, GEM and HOMER), is more successful in dealing with such a dataset and calls a

210     total of 12,452 peaks. We therefore sought to investigate the efficacy of utilizing the union peak set

211     for this poorly enriched SALL4 ChIP-Seq, which enables us to sacrifice specificity for a gain in

212     sensitivity across the dataset, ie, we accept all peaks including those called by only a single peak

213     caller which carry less confidence but provide higher sensitivity (**Supplemental Table 5**).

214

215     To test its validity, we compared the SALL4 ChIP-Seq union peak set with a SALL4 Cut&Run

216     dataset recently published[33]. Cut&Run is a technique which utilizes antibody-targeting and

217     micrococcal nuclease digestion to map global DNA binding sites[34]. It is an analogous but independent

218     technique to ChIP-Seq thus providing an independent dataset for comparison and validation. To

219    ensure the peaks called in our ChIP-Seq were likely binding targets of SALL4, we first directly

220    compared the union peak set to the SALL4 Cut&Run dataset. Reassuringly, the union peak set

221    showed a 36% overlap with peaks identified in at least 2 of the Cut&Run replicates (3 biological

222    replicates total) (**Figure 2c**). Using individual peak callers, overlap percentages ranging from 25-56%

223    were observed with fewer peaks called (**Supplemental Figure 2a**). Furthermore, each caller reports a

224    different sub-set of targets with little overlap between them (**Figure 2b**). However, by considering the

225    union peak set, we can gain a more complete overview of the binding landscape without significantly

226    affecting average sensitivity, by allowing us to circumvent the poor performance of individual peak

227    callers for the dataset in question and call "missed" peaks.

228

229    Next, we wanted to confirm that the called union peaks show our recently identified human SALL4

230    DNA binding motif[33]. To investigate this, we performed a directed motif search wherein we searched

231    every peak in the union peak set for the human SALL4 DNA binding motif. This showed that 55% of

232    the union peak set contained at least one instance of our identified SALL4 motif (**Figure 2d**), a result

233    comparable to using an individual peak caller alone (**Supplemental Table 1, 6**). To ensure we have

234    not biased the motif search, we performed a *de novo* motif search on the union peak set using both the

235    MEME-Suite[22-24] (which utilize the algorithms STREME, CentriMo and MEME-ChIP) and HOMER

236    which were both able to call the human SALL4 DNA binding motif as the second and third top

237    candidate motif hits respectively (**Figure 2e, f**). According to CentriMo, the STREME identified

238    motif is centrally enriched in individual peaks in the union peak set (**Figure 2g**), an expected

239    observation for true binding motif sequences. Furthermore, MEME-ChIP itself, called the same motif

240    as the third candidate peak with the second candidate motif result also being an AT rich motif with

241    near identical sequence (**Supplemental Figure 2b**). We therefore concluded that despite the

242    additional number of peaks called by taking the union peak set, the SALL4 DNA-binding motif

243    signature is still present across all called peaks and is identifiable using multiple algorithms as a top

244    three candidate motif.

245

246    As further validation to ensure the union peak set identified valid targets of SALL4, a GO analysis

247    was performed and compared with the results of previous findings[33]. We previously reported that

248    SALL4 knock-down resulted in a significant increase in the number of up-regulated genes in the

249    "*transcriptional regulator activity*" (GO:0140110) pathway, results which were validated by

250    comparing bulk RNA-Seq and Cut&Run results[33], and thus confirming pathway members as *bona fide*

251    SALL4 targets. Consistently, the GO analysis on the union peak set identified the same pathway as a

252    top 20 enriched pathway (**Supplemental Table 7**), with more significantly enriched terms pointing to

253    SALL4 being a DNA-binding protein; a well-established function of SALL4[31,32] (**Supplemental**

254    **Table 7**). Therefore, despite utilizing the union peak set which sacrificed a degree of specificity, the

255    peak set was still valid in detecting accurate biological functions of SALL4.

256

257    The final validation to ensure the union dataset identified valid targets of SALL4 was to overlap the

258    union peaks gene list with the SALL4 knock-down bulk RNA-Seq previously published[33], and

259    compare the overlap targets of the union peak set with the overlapping targets of the Cut&Run

260    dataset. We previously reported that 2,695 genes were found significantly differentially expressed on

261    SALL4 knock-down, 430 of which has a corresponding annotated SALL4 Cut&Run peak. Using the

262    union peak set, we observed an overlap of 451 genes (**Supplemental Table 8**) with the SALL4

263    knock-down dataset, of which, 198 gene targets were found in common between the Cut&Run and

264    ChIP-Seq gene sets (**Supplemental Table 8**). This finding combined with the observed overlap

265    between the union peak set and the Cut&Run replicates suggests that there are SALL4 binding targets

266    that were detected by the ChIP-Seq that were not detected by the Cut&Run and vice versa. However,

267    both the Cut&Run and the union peak list derived from the SALL4 ChIP-Seq appear to still be calling

268    valid SALL4 target genes with significant overlaps between the two datasets observed.

269

270    Taken together, the results obtained show that despite this SALL4 ChIP-Seq showing poor

271    enrichment with few peaks called by two of the four peak callers, there was still valid data within the

272    dataset that can be extracted, used, and validated by independent approaches[33]. By considering the

273    union peak set generated by ChIP-AP, one can opt to marginally sacrifice specificity for a significant

274    gain in sensitivity across the dataset and confirm the presence of peaks identified or validated using

275    different methodologies should the characteristics of the dataset prove to be less than favourable.

276    Whereas previous analyses using a single peak caller would produce sub-optimal results, by relying

277    on multiple peak callers, as ChIP-AP does, sub-optimal datasets can be salvaged and still report valid

278    findings.

279

## 280    ChIP-AP Functionality and Characteristics
### 281    ChIP-AP Modularity for Advanced Users

282    Many programming languages are based on the programming paradigm of Object-Oriented

283    Programming (OOP), wherein individual components of the program resemble "reusable objects"

284    with defined input and output parameters (**Figure 3a**). This compartmentalization allows the

285    programmer to assemble these "objects" in any manner to accomplish the task at hand provided the

286    requisite parameters are met for individual objects. In the same spirit as OOP, ChIP-AP has been

287    designed to be "object-oriented" in nature (**Figure 3b**).

288

289    To instantiate a ChIP-AP run, all input arguments are passed through the command line or the

290    graphical interfaces. What is essential for a run is the location of the input sequencing files (raw fastq)

291    and a settings table for customization of pipeline constituent programs (discussed in the following

292    section) (**Figure 3b**). ChIP-AP then proceeds to then construct a folder hierarchy and places within

293    each folder the corresponding sub-script for that stage of analysis. Each ChIP-AP sub-script script is

294    in essence an instantiated object with defined input and output parameters passing files sequentially

295    from one folder to the next for processing and analysis. Should a user wish to add to or remove an

296    aspect of the pipeline, one simply needs to be mindful of the adjoining objects input/output

297    characteristics. ChIP-AP therefore provides an analysis platform wherein individual analysis steps can

298    be modularly swapped with equivalent steps, provided they have identical input and output

299    characteristics, without requiring additional changes to the flow of analysis or code. This

300    compartmentalization of analytical steps enables ChIP-AP to be exponentially customizable to

301    differing scenarios if the user is proficient enough to code the equivalent analysis step required. The

302    ChIP-AP documentation on GitHub accurately outlines all the analysis steps and documents the input

303    and output behaviours of each sub-script, this is in addition to a comprehensively commented master

304    script outlining the same information in code.

305

306    Constituent Program Customization and Analysis Reproducibility Through the Settings Table
307    The lack of result reproducibility in science is a major and on-going issue[35]. The field has continued to

308    change and adapt to this problem with journals enforcing stricter reporting of materials and methods

309    in an attempt to curtail such issues. Unfortunately, bioinformatics methods reporting is an area of

310    scientific research where significant work is still required. Reporting of ChIP-Seq analyses in

311    publications consistently lacks necessary details with many authors omitting key program

312    modification parameters or even neglecting to mention key analysis steps entirely. We have therefore

313    attempted to address this issue by ensuring ChIP-AP analyses are reproducible through an accurate

314    and consistent means of reporting.

315

316    A key design aspect of ChIP-AP was to require the provision of a Settings Table (ST). If no table is

317    provided, ChIP-AP will use a pre-defined default-ST (DST; **Table 2**). The ST lists all the programs

318    used in the ChIP-AP run and all the necessary optional program arguments entered for that particular

319    run. It is therefore a listing of all non-hard-coded program modification parameters/flags used for a

320    particular analysis. For ChIP-AP to reproduce any analysis, it simply needs the raw sequencing fastq

321    files and the ST used. *We consider the dissemination of the information contained in the ST as both*

322    *vital and essential, along with results obtained*. The ST can be included as a supplemental table in a

323    manuscript or can be included as a processed data file when submitting data to an upload repository

324    like GEO. In either case, the information of this file *must* be presented when publishing data to ensure

325    analysis reproducibility in a format that is both consistent and convenient. Of note also, whether the

326    user provided a ST as input or the default-ST was used, a copy of the table will be found in the output

327    folder to ensure all required program modification parameters are provided accompanying the final

328    result.

329

## ChIP-AP User Interfaces for Biologists

330

331    There is an ever-increasing need to make dry-lab analyses accessible to wet-lab biologists wishing to

332    investigate and interrogate data themselves without having to collaborate with (or wait for) a

333    bioinformatician. This is straightforward if a single program is required for an analysis like GraphPad

334    Prism or SAS. However, longer or more comprehensive analyses and workflows would typically

335    require a degree of coding to work. It is in this domain that ChIP-Seq analysis resides as it requires

336    the utilization of multiple programs, each feeding into each other to perform a coherent analysis. In an

337    attempt to address this demand, platforms such as Galaxy, or licensed software such as the Partek

338    Genomics Suite, have been developed to add graphical user interface (GUI) elements to analyses to

339    make higher-level analyses more accessible to researchers with no coding background. These

340    platforms though, particularly for ChIP-Seq analyses, utilize only a single peak caller and can offer

341    limited customization of program parameters in certain scenarios. As discussed, this can result in

342    incomplete analysis of the bound landscape owing to algorithmic limitations and biases, issues

343    ChIP-AP was designed to address. It was therefore necessary for ChIP-AP to incorporate its own GUI

344    to aide users in completing their required analyses and thus enable researchers with no coding

345    experience to perform independent analyses.

346

347    To address the breadth of computer proficiencies seen in the wet-lab scientific community, we

348    implemented two GUI's, the choice of which to use will depend on a user's proficiency with

349    ChIP-AP. Through the guided step-by-step tutorials found on our Github, users can install ChIP-AP

350    on any modern operating system, including Windows 10, and run the GUI of their choosing. The

351    GitHub repository lists the system hardware requirements to run ChIP-AP, but many modern laptops

352     and computers commonly purchased for research are capable of analysing data locally, without

353     needing dedicated server hardware.

354

355     For users unfamiliar with the command line, we have implemented the Wizard interface (**Figure 3c**),

356     inspired by installation wizards from the Windows 95/98 era of computing. The ChIP-AP wizard will

357     guide users through the analysis configuration by means of a series of panels each asking a single

358     question about the input data. On completion, users will have provided all the necessary information

359     required for a ChIP-AP run and can start their analysis directly from the wizard. This GUI

360     implementation was designed to not overwhelm users with multiple questions simultaneously asking

361     for input, but rather asks for data in a more guided approach.

362

363     For users familiar with the input requirements of ChIP-AP, we have implemented the Dashboard

364     interface (**Figure 3d**). The dashboard asks the same questions as the wizard but in a single panel,

365     enabling users to input the required data more quickly (**Figure 3d** – Data Input). Once all the required

366     information is input, as with the wizard, users can run ChIP-AP directly from the interface. In stark

367     difference to the wizard though, the dashboard interface contains a command line translation window

368     at the bottom of the interface (**Figure 3d** – Command Line Translation). As users enter data in the

369     GUI elements, the command line translation window will automatically update to accommodate the

370     additional/changed inputs. This enables researchers to gradually draw connections between translating

371     static GUI elements into command line arguments and flags to modulate and control program

372     behaviour. Such an implementation will aide some researchers more comfortably and confidently

373     transition from GUI to command line usage of ChIP-AP, and hopefully, beyond for their research.

374

375     Finally, independent of whether a user opts to use the wizard or dashboard interface, users will be

376     prompted to either use the DST or choose to upload their own ST. As discussed, the functionality and

377    reproducibility provided by the DST/ST is essential for ChIP-AP reproducibility, thus enabling a GUI

378    utilizing researcher to reproduce an analysis performed and customized by a bioinformatician.

379

380    Conclusion

381    ChIP-Seq is a well-established experimental protocol for profiling DNA-interacting proteins. In the

382    computational space, many software tools have been developed with over 50 peak callers being

383    published to date. Despite the abundance of available peak callers however, benchmarking studies

384    have consistently shown poor overlap between peak-sets from different peak callers. This is because

385    every caller has distinct selectivity and specificity characteristics which are often not additive and

386    seldom completely overlap with other peak callers in many scenarios. Additionally, it has been

387    extensively shown that the performance of a peak caller is subject to the read characteristics and read

388    distributions of the dataset in question. An individual peak caller can outperform other callers in

389    certain datasets but will perform poorly in alternate datasets. Therefore, with the heterogeneity

390    observed in experimental samples arising from profiling different DNA-binding proteins each profiled

391    with differing immunoprecipitation and library preparation protocols, reliance on a single peak caller

392    is unlikely to yield the most reliable, consistent or comprehensive results.

393

394    To circumvent the limitations and biases of individual peak callers, we rationalized that integrating

395    the results of multiple peak callers would yield improved peak calling consistency, confidence and

396    more comprehensively assess the binding landscape without requiring additional wet-lab

397    observations. As such, we developed the integrated ChIP-Seq analysis pipeline, ChIP-AP, which takes

398    design decisions from established workflows such as those utilized in consortia projects like

399    ENCODE[18]. ChIP-AP has been coded from the ground-up to be as simple to use as possible for users

400    inexperienced with the command line by providing two GUI's for use, the wizard or dashboard.

401    ChIP-AP still however remains exponentially customizable for advanced of users by facilitating

402    fine-grained customization of constituent programs through the ST, or, through its provision of

403    customizable modular framework that enables swapping of analysis stages to tailor ChIP-AP for

404    custom workflows. While ChIP-AP has been designed and written specifically for ChIP-Seq analysis,

405    the framework and design principles on which its coded, facilitate its adaptation and use for other

406    existing (ATAC-Seq[36,37], RIP-Seq[38], Cut&Run[34]) and future emerging technologies. Should a new,

407    peak caller or analytical tool be developed, minimal changes are required to add an additional step in

408    the pipeline to accommodate the inclusion of said tool. This allows ChIP-AP to be easily modified to

409    work with emerging techniques and any tools that will be specifically developed for such a technique.

410    ChIP-AP can therefore be expanded or enhanced to suit future applications and uses with necessary

411    program arguments being passed through the settings table for each ChIP-AP run. To the best of our

412    research, we have yet to find an integrated software solution currently available that utilizes multiple

413    peak callers other than ChIP-AP.

414

415    ChIP-AP as presented here though, allows users to sub-set the binding landscape in a manner that is

416    best suited to address their biological research question, while allowing users to switch between

417    differing sub-sets depending on the question at hand. By utilizing the consensus peak set, binding

418    motif accuracy can be significantly increased by restricting the motif search space to only the most

419    confident peaks. This also has improved outcomes when attempting down-stream GO analysis

420    wherein more targeted and biologically significant terms can be reported. In contrast though, if the

421    profiled dataset has unfavourable characteristics such as poor enrichment or shows high signal:noise,

422    the union peak set can potentially yield improved results and allow users to marginally sacrifice

423    specificity for a potentially significant increase in sensitivity across the binding landscape. In between

424    these two extremes of data sub-sets is a gradient of sensitivity thresholds that can be selected

425    depending on the biological question and the presence of additional, supportive data from independent

426    techniques and methodologies. By reporting such an integrated analysis, ChIP-AP enables the end

427    user to focus on the biological question at hand by providing a comprehensive protein binding profile

428    without needing data re-analysis. ChIP-AP can therefore provide both substantial improvements to

429    peak capturing and analysis reliability from a single integrated and comprehensive analysis.

430

431    ## Materials and Methods
432    ## ChIP-AP Constituent Programs
433    ChIP-AP is an integrated pipeline that brings multiple command line programs together into a single,

434    seamless and easy to use pipeline. At time of publication, these include FastQC[39], Clumpify and

435    BBDuk from the BBMap Suite[40], Trimmomatic[41], BWA[42], Samtools[43], deepTools[44], MACS2[45],

436    GEM[16], SICER2[46], HOMER[15] and Genrich[13]. If using ChIP-AP, please cite all constituent tools as

437    well. It is best to refer to the GitHub repository for the latest citation list which would include any

438    additional tools incorporated into ChIP-AP since publication.

439

440    ## SNU-398 Culturing
441    SNU-398 cell line was obtained from the American Type Culture Collection (ATCC). The cells were

442    maintained in RPMI medium supplemented with 10% foetal bovine serum (FBS) at 37°C in a

443    humidified atmosphere of 5% CO2 as recommended by ATCC.

444

445    ## SALL4 ChIP-Seq Preparation and Sequencing
446    20 million SNU-398 cells were cross-linked with 1% formaldehyde for 10 minutes at room

447    temperature. The reaction was terminated by adding 2M glycine to a final concentration of 125mM.

448    Cells were then washed with 1×PBS and resuspended in 1mL of cell lysis buffer (20mM Tris pH8.0,

449    85mM KCl, 0.5% nonidet P-40, protease inhibitor). After 10 minutes of incubation on ice, cells were

450    spun down and cell pellet was resuspended in another 1mL of cell lysis buffer. After another 5

451    minutes of incubation on ice, cells were spun down and cell pellet was resuspended in 1mL of nuclear

452    lysis buffer (10mM Tris-HCl pH7.5, 1% nonidet P-40, 0.5% sodium deoxycholate, 0.1% SDS,

453    protease inhibitor). After 10 minutes of incubation on ice, chromatin was sheared to 500bp.

454    Antibody-protein A/G Dynabead conjugate was prepared by adding 0.75μg of SALL4 rabbit

455    monoclonal antibody (Cell Signaling Technology #8459) to pre-washed 50μL of protein A/G

456    Dynabeads (Life Techonlogies) with one hour incubation at 4°C with rotation. Sheared chromatin was

457    then added to antibody-protein A/G conjugate and incubated overnight at 4°C with rotation. After

458    overnight incubation, the beads were washed sequentially with the following buffers: twice with

459    RIPA/500mM NaCl buffer (0.1% deoxycholate, 0.1% SDS, 1% Triton X-100, 500mM NaCl, 1mM

460    EDTA, 20mM Tris-HCl pH8.1), twice with LiCl buffer (0.25M LiCl, 1% nonidet P-40, 1% sodium

461    deoxycholate, 1mM EDTA, 10mM Tris-HCl pH8.1), twice with TE buffer (10mM Tris-HCl pH8.0,

462    1mM EDTA pH8.0). Protein complexes were reverse cross-linked with 50µL of ChIP Elution Buffer

463    (10mM Tris-HCl pH8.0, 5mM EDTA, 300mM NaCl, 0.1% SDS) and 8µL of Reverse Crosslink Mix

464    (250mM Tris-HCl pH6.5, 1.25M NaCl, 62.5mM EDTA, 5mg/mL proteinase K, 62.5µg/mL RNase A)

465    at 65°C for 5 hours. Reverse cross-linked DNA was cleaned up using SPRI beads (Beckman Coulter)

466    and eluted in 10mM Tris-HCl pH 8.0. To generate libraries for deep sequencing, the eluted DNA was

467    end-repaired using End-It DNA End-Repair Kit (Epicenter #ER0720) and A-tailing was then carried

468    using Klenow (3'-5' exo-) enzyme (New England Biolabs). Illumina sequencing adaptors were ligated

469    to the DNA fragments and adaptor-ligated DNA fragments were enriched with 14 cycles of PCR.

470    DNA libraries were gel purified and analyzed on Bioanalyzer (Agilent) for their size distribution.

471    Libraries were sequenced on Illumina HiSeq 2500 sequencer with single-end 35bp settings.  The

472    sequencing and processed files have been uploaded to GEO with Accession number xxxx (reviewer

473    access token xxxx).

474

475    SALL4 ChIP-Seq Analysis and Comparisons
476    The generated SALL4 ChIP-Seq was processed with ChIP-AP (v4.1) using a hg38 genome. The

477    settings table used for analysis is found below.

| Program | Argument |
| --- | --- |
| fastqc1 | -q |
| clumpify | dedupe spany addcount |
| bbduk | ktrim=l hdist=2 |
| trimmomatic | LEADING:20 SLIDINGWINDOW:4:20 TRAILING:20 MINLEN:20 |
| fastqc2 | -q |
| bwa_mem | |
| samtools_view | |
| plotfingerprint | -bs 50 --centerReads –ignoreDuplicates |
| fastqcs3 | -q |
| macs2_callpeak | |
| gem | -Xmx30G --k_min 8 --k_max 12 |
| sicer2 | |
| homer_findPeaks | |
| genrich | --adjustp -v |

| homer_mergePeaks | |
|---|---|
| homer_annotatePeaks | |
| fold_change_calculator | --normfactor uniquely_mapped |

478

479     For all analyses, the union peak set was utilized. The fingerprint plot (Figure 2a) was generated as part

480     of the ChIP-AP run with the flags outlined in the settings table. The upset plot (Figure 2b) was

481     generated by taking the "venn.txt" data from the ChIP-AP run output (folder 21_peaks_merging) and

482     plotting it in R[47] (v4.0.3) with the UpSetR[48] (v1.4.0) package.

483     For comparisons with the Cut&Run data, the Cut&Run data was processed as outlined previously[33]

484     and is available from GEO, Accession GSE136332. To overlap the Cut&Run replicates, HOMER's

485     mergePeaks was used with flags "-d 1500." Next, the Cut&Run peaks identified in at least 2 replicates

486     were combined into a single list and compared to the SALL4 ChIP-Seq union peak set using

487     HOMER's mergePeaks with flag "-d 2000", this provided the list of overlapping regions, the number

488     of which was plotted in R[47] (v4.0.3) with the VennDiagram[49] (v1.6.20) package.

489     For the directed motif search within the SALL4 ChIP-Seq union peak set, HOMER's[15]

490     findMotifsGenome function was used with flags "-find sall4_weighted_motif.motif." For the

491     HOMER *de novo* motif search, HOMER's findMotifsGenome function was used with flags "hg38 -

492     size given -mask." For the MEME-ChIP (and sub-program[22-24]) motif search, first the union peak list

493     was processed with HOMER's findMotifsGenome function with flag "-dumpFasta" to extract the

494     central 200bp sequences of each peak.  HOMER also generated an equivalent set of background

495     sequences with comparable GC content to be used. Next, MEME-ChIP was run with flags "-neg

496     background.fa -meme-nmotifs 25 union_peaks.fa." Motif logo files were generated using R[47] (v4.0.3)

497     and the seqLogo[50] (v1.52.0) package.

498     The gene ontology analysis of the SALL4 ChIP-Seq dataset performed was part of the ChIP-AP run

499     using the flag "-goann" which utilizes HOMER to perform the analysis following peak annotation. To

500     compare with the processed SALL4 knock-down results published[33], we started from supplemental

501     tables 4 and 5 from the publication. Next, we overlayed the reported gene names from the SALL4

502     ChIP-Seq union peak set to those gene lists to determine overlapping gene names.

503

## Encode Datasets Utilized and Processing

505     A number of ENCODE datasets were downloaded and utilized for our analysis. The table below lists

506     all the downloaded experiment ID's used. Data was downloaded from ENCODE March 2021.

| Cell Line | Transcription Factor | ChIP Experiment ID's | Control Experiment ID's |
|---|---|---|---|
| GM12878 | MAX | ENCFF000VXY ENCFF000VYA | ENCFF000VWF ENCFF000VWH |
|  | SPI1 | ENCFF000VXY ENCFF000VYA |  |
| HepG2 | ZBTB33 | ENCFF000PSP ENCFF000PSW | ENCFF000POC ENCFF000POH |
|  | CEBPB | ENCFF000XQM ENCFF000XQN |  |
| K562 | MAFF | ENCFF000YSQ ENCFF000YSS | ENCFF002EFF ENCFF002EFD |
|  | JUN | ENCFF000YJJ ENCFF000YJL |  |
|  | GATA1 | ENCFF000YND ENCFF000YNF |  |
|  | MEIS2 | R1: ENCFF002EIU       ENCFF002EIW R2: ENCFF002EIV       ENCFF002EIX | R1: ENCFF002EFF       ENCFF002EFD R2: ENCFF002EFH       ENCFF002EFA |
|  | RUNX1 | R1: ENCFF002DOZ       ENCFF002EGD R2: ENCFF002EGE       ENCFF002DPH |  |
|  | ATF4 | R1: ENCFF081USS       ENCFF565KLI R2: ENCFF069VNL       ENCFF682IGK |  |

507

508     All ENCODE datasets were processed with ChIP-AP (v4.1) with the following settings table

| Program | Argument |
|---|---|
| fastqc1 | -q |
| clumpify | dedupe spany addcount qout=33 fixjunk |
| bbduk | ktrim=l hdist=2 |
| trimmomatic | LEADING:20 SLIDINGWINDOW:4:20 TRAILING:20 MINLEN:20 |
| fastqc2 | -q |
| bwa_mem |  |
| samtools_view |  |
| plotfingerprint | -bs 50 --centerReads –ignoreDuplicates |
| fastqcs3 | -q |
| macs2_callpeak |  |
| gem | -Xmx30G --k_min 8 --k_max 12 |
| sicer2 |  |
| homer_findPeaks |  |
| genrich | --adjustp -v |

| | |
|---|---|
| homer_mergePeaks | |
| homer_annotatePeaks | |
| fold_change_calculator | --normfactor uniquely_mapped |

509     For each transcription factor, the corresponding JASPAR binding motif for the cell line in question

510     was downloaded from MethMotif[19] and manually converted to HOMER motif format. For the

511     directed motif searches, HOMER's findMotifsGenome was utilized with flags "hg38 -find

512     binding_motif.motif." For HOMER *de novo* motif discovery, the findMotifsGenome program was

513     used with flags "hg38 -size given -mask -dumpFasta." This ran the motif discovery while also giving

514     the necessary fasta sequence files (target.fa and background.fa) required to run the MEME-Suite. The

515     MEME *de novo* motif discovery was run with flags "-neg background.fa -meme-nmotifs 25 target.fa."

516     Motif logo files were generated using $R^{47}$ (v4.0.3) and the seqLogo[50] (v1.52.0) package. The gene

517     ontology results were generated as part of HOMER's annotatePeaks function for the required peak

518     sets. HOMER annotatePeaks was utilized with a known motif provided with -m flag  to include the

519     distances from all starting coordinate motif instances in each peak to their respective peak starting

520     coordinate. A custom script was utilized to extract the distances from every peak's weighted peak

521     center coordinate to the midpoint coordinate of the motif instance closest to the weighted peak center.

522     The density plots representing this data were generated using $R^{47}$ and the ggplot2[51]. Peak-Motif

523     percentages were plotted using Graphpad Prism v9.1.0.

524

525 # Figure and Table Legends

526 ## Figure 1 – Consensus peak set improves detected motif accuracy

527 a) Peak-Motif percentage (number of peaks with binding motif) for all 10 TF's profiled as identified

528 for the MACS2 and the consensus peak sets. b) The motif position-bias for CEBPB, JUN, SPI1 and

529 ZBTB33 for the consensus peak set and the individual peak callers. The position-bias is a measure of

530 how far the identified motif sits away from the weighted peak center. c) The CEBPB *de novo* motif

531 discovery results as reported by HOMER for the consensus peak set. The line above the peaks

532 delineates position of the binding motif. d) The CEBPB *de novo* motif discovery results as reported by

533 HOMER for the MACS2 peak set. The line above the peaks delineates position of the binding motif.

534 e) The MAFF *de novo* motif discovery results as reported by HOMER for the consensus peak set. The

535 line above the peaks delineates position of the binding motif. f) The MAFF *de novo* motif discovery

536 results as reported by HOMER for the MACS2 peak set. The line above the peaks delineates the

537 position of the binding motif.

538

539 ## Figure 2 – Union peak set comprehensiveness and accuracy

540 a) Fingerprint plot for aligned sequence files for samples. Negligible separation between the SALL4

541 and control curves indicates poor enrichment in the SALL4 samples. b) Upset plot describing the

542 distribution of peaks observed by each peak caller. The left histogram represents the total number of

543 called peaks per caller. The top histograms represent the size of the sub-sets in question. The

544 connected circles represent highlighted overlap. c) Venn diagram showing the overlapping number of

545 peaks between the SALL4 union ChIP-Seq dataset and the Cut&Run dataset. d) The motif sequence

546 used for the directed motif search in the SALL4 ChIP-Seq union set, which was found in 55.2% of the

547 union set. e) The STREME *de novo* motif search for the SALL4 union peak set identified the AT-rich

548 binding motif as the 2nd result. f) The HOMER de novo motif search for the SALL4 union peak set

549 identified the AT-rich binding motif as the 3rd result. g) The STREME identified motif (shown in d)

550 was found centrally enriched in the union peak set as compared to background sequences.

551

552 ## Figure 3 – Object Oriented Nature of ChIP-AP

553 a) In OOP, an abstract "object" is defined as a segment of code that accepts defined inputs, processes

554 the data, and outputs the data in a defined manner. Objects can then be combined in any manner to

555 produce desired output. b) ChIP-AP was designed to be "object-oriented" in nature with each stage of

556 analysis in a folder (01, 02…) having defined input/output characteristics. c) The ChIP-AP wizard

557 interface guides users through a series of windows, each asking for a single piece of input, till all

558 required information is gathered for a ChIP-AP run. d) The dashboard interface can be separated into

559 2 regions, the data input and command line translations segments. In the data input section, all the

560    required data for a ChIP-AP run is input from a single interface. The command line translation

561    window at the bottom dynamically changes as input is entered in the data input section, translating

562    static GUI elements into the necessary command line arguments/flags enabling users to view how

563    ChIP-AP's input is modified based on the provided input.

564

## Supplemental Figure 1

566    a) The motif position-bias for ATF4, GATA1, MAFF, MAX, MEIS2 and RUNX1 for the consensus

567    peak set and the individual peak callers. The position-bias is a measure of how far the identified motif

568    sits away from the weighted peak center. b) The CEBPB *de novo* motif discovery results as reported

569    by the MEME-Suite for the consensus peak set. Above each p-value is which sub-program of MEME

570    called the said motif. c) The CEBPB *de novo* motif discovery results as reported by the MEME-Suite

571    for the MACS2 peak set. Above each p-value is which sub-program of MEME called said motif. d)

572    The MAFF *de novo* motif discovery results as reported by the MEME-Suite for the consensus peak

573    set. Above each p-value is which sub-program of MEME called said motif. The $6^{th}$ result shows a

574    characteristic heterodimer binding profile for MAFF. e) The MAFF *de novo* motif discovery results as

575    reported by the MEME-Suite for the MACS2 peak set. Above each p-value is which sub-program of

576    MEME called said motif.

577

## Supplemental Figure 2

579    a) Venn diagrams highlighting the degree of overlap between each individual callers peak-set and the

580    Cut&Run peak set, and the relatively few peaks called by each individual peak caller. b) The

581    MEME-ChIP results highlighting showing the correct AT-rich binding motif for SALL4 is the $3^{rd}$

582    called motif hit. The $2^{nd}$ motif hit also is an AT-rich motif with near identical sequence.

583

## Table 1 – Top 20 RUNX1 GO Terms

585    The top 20 RUNX1 GO terms returned for the consensus peak set (left) and the MACS2 peak set

586    (right). The consensus peak set GO returned GO terms are more directly relatable to defined RUNX1

587    functions as compared to the MACS2 results.

588

## Table 2 – Default Settings Table for ChIP-AP

590    The default program settings table used by ChIP-AP if no user provided settings table is provided.

591    The left column lists the constituent programs of ChIP-AP with their optional modification

592    parameters/flags found in the right column.

593

594 ## Supplemental Table 1
595 Table listing and overview of all the profiled TF's, their TF family and cell line of origin. The table

596 also lists the total number of peaks found in the MACS2 and consensus peak sets, along with the

597 peak-motif percentages for each set.

598

599 ## Supplemental Table 2
600 A listing of the z-tests performed testing the position-bias distributions of the consensus peak set

601 compared to each individual peak caller. Significant differences are highlighted in green. Cells

602 highlighted yellow indicate values approaching significance.

603

604 ## Supplemental Table 3
605 All the consensus peak GO results for all 10 TF's (1 sheet per TF).

606

607 ## Supplemental Table 4
608 All the MACS2 peak GO results for all 10 TF's (1 sheet per TF).

609

610 ## Supplemental Table 5
611 The union peak-list for the SALL4 ChIP-seq

612

613 ## Supplemental Table 6
614 The peak-motif percentages for each individual peak callers results in the SALL4 ChIP-Seq dataset.

615

616 ## Supplemental Table 7
617 The GO results for the SALL4 union peak-set.

618

619 ## Supplemental Table 8
620 The differentially expressed genes from the SALL4 knock-down RNA-Seq experiment found to

621 contain at least 1 peak in the SALL4 union peak-set.

622

## Acknowledgements

## Contributions

J.S. and M.A.B designed the package. J.S. was the lead programmer. J.S. and M.A.B tested, optimized and debugged ChIP-AP. D.E.T and KJ.Y. performed the SALL4 ChIP-Seq. J.S., D.G.T and M.A.B interpreted results and wrote the manuscript. M.A.B conceived and directed the project.

# References

1    Collas, P. The current state of chromatin immunoprecipitation. *Mol Biotechnol* **45**, 87-100, doi:10.1007/s12033-009-9239-8 (2010).

2    Johnson, D. S., Mortazavi, A., Myers, R. M. & Wold, B. Genome-wide mapping of in vivo protein-DNA interactions. *Science* **316**, 1497-1502, doi:10.1126/science.1141319 (2007).

3    Bernstein, B. E. *et al.* Genomic maps and comparative analysis of histone modifications in human and mouse. *Cell* **120**, 169-181, doi:10.1016/j.cell.2005.01.001 (2005).

4    Park, P. J. ChIP-seq: advantages and challenges of a maturing technology. *Nat Rev Genet* **10**, 669-680, doi:10.1038/nrg2641 (2009).

5    Barski, A. *et al.* High-resolution profiling of histone methylations in the human genome. *Cell* **129**, 823-837, doi:10.1016/j.cell.2007.05.009 (2007).

6    Mikkelsen, T. S. *et al.* Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature* **448**, 553-560, doi:10.1038/nature06008 (2007).

7    Chen, Y. *et al.* Systematic evaluation of factors influencing ChIP-seq fidelity. *Nat Methods* **9**, 609-614, doi:10.1038/nmeth.1985 (2012).

8    Kurzawa, N., Eils, R., Steinhauser, S. & Herrmann, C. A comprehensive comparison of tools for differential ChIP-seq analysis. *Briefings in Bioinformatics* **17**, 953-966, doi:10.1093/bib/bbv110 (2016).

9    Laajala, T. D. *et al.* A practical comparison of methods for detecting transcription factor binding sites in ChIP-seq experiments. *BMC Genomics* **10**, 618, doi:10.1186/1471-2164-10-618 (2009).

10   Wilbanks, E. G. & Facciotti, M. T. Evaluation of algorithm performance in ChIP-seq peak detection. *PLoS One* **5**, e11471, doi:10.1371/journal.pone.0011471 (2010).

11   Koohy, H., Down, T. A., Spivakov, M. & Hubbard, T. A comparison of peak callers used for DNase-Seq data. *PLoS One* **9**, e96303, doi:10.1371/journal.pone.0096303 (2014).

12   Jeon, H., Lee, H., Kang, B., Jang, I. & Roh, T. Y. Comparative analysis of commonly used peak calling programs for ChIP-Seq analysis. *Genomics Inform* **18**, e42, doi:10.5808/GI.2020.18.4.e42 (2020).

13   Gaspar, J. *Genrich: detecting sites of genomic enrichment*, <https://github.com/jsh58/Genrich> (

14   Liu, T. Use model-based Analysis of ChIP-Seq (MACS) to analyze short reads generated by sequencing protein-DNA interactions in embryonic stem cells. *Methods Mol Biol* **1150**, 81-95, doi:10.1007/978-1-4939-0512-6_4 (2014).

15   Heinz, S. *et al.* Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol Cell* **38**, 576-589, doi:10.1016/j.molcel.2010.05.004 (2010).

16   Guo, Y., Mahony, S. & Gifford, D. K. High resolution genome wide binding event finding and motif discovery reveals transcription factor spatial binding constraints. *PLoS Comput Biol* **8**, e1002638, doi:10.1371/journal.pcbi.1002638 (2012).

17   Xu, S., Grullon, S., Ge, K. & Peng, W. Spatial clustering for identification of ChIP-enriched regions (SICER) to map regions of histone methylation patterns in embryonic stem cells. *Methods Mol Biol* **1150**, 97-111, doi:10.1007/978-1-4939-0512-6_5 (2014).

18   Consortium, E. P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74, doi:10.1038/nature11247 (2012).

19   Xuan Lin, Q. X. *et al.* MethMotif: an integrative cell specific database of transcription factor binding motifs coupled with DNA methylation profiles. *Nucleic Acids Res* **47**, D145-D154, doi:10.1093/nar/gky1005 (2019).

20   Khan, A. *et al.* JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Res* **46**, D260-D266, doi:10.1093/nar/gkx1126 (2018).

21   Lin, Q. X. X., Thieffry, D., Jha, S. & Benoukraf, T. TFregulomeR reveals transcription factors' context-specific features and functions. *Nucleic Acids Res* **48**, e10, doi:10.1093/nar/gkz1088 (2020).
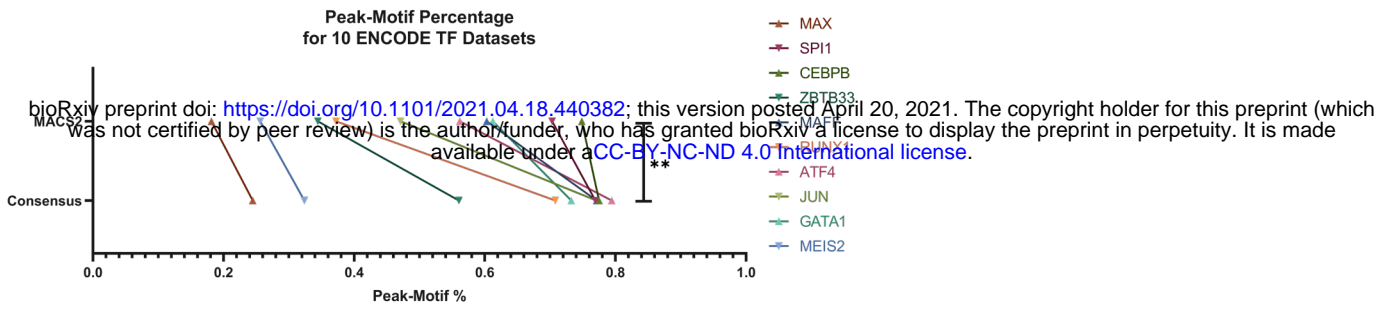
688    22    Bailey, T. L. STREME: Accurate and versatile sequence motif discovery. *bioRxiv*,
689          2020.2011.2023.394619, doi:10.1101/2020.11.23.394619 (2020).
690    23    Bailey, T. L. & Machanick, P. Inferring direct DNA binding from ChIP-seq. *Nucleic Acids*
691          *Res* **40**, e128, doi:10.1093/nar/gks433 (2012).
692    24    Machanick, P. & Bailey, T. L. MEME-ChIP: motif analysis of large DNA datasets.
693          *Bioinformatics* **27**, 1696-1697, doi:10.1093/bioinformatics/btr189 (2011).
694    25    Fuka, G. *et al.* Silencing of ETV6/RUNX1 abrogates PI3K/AKT/mTOR signaling and
695          impairs reconstitution of leukemia in xenografts. *Leukemia* **26**, 927-933,
696          doi:10.1038/leu.2011.322 (2012).
697    26    Imperato, M. R., Cauchy, P., Obier, N. & Bonifer, C. J. I. J. o. H. The RUNX1–PU.1 axis in
698          the control of hematopoiesis.  **101**, 319-329, doi:10.1007/s12185-015-1762-8 (2015).
699    27    Lam, K. & Zhang, D. E. RUNX1 and RUNX1-ETO: roles in hematopoiesis and
700          leukemogenesis. *Front Biosci (Landmark Ed)* **17**, 1120-1139 (2012).
701    28    Pencovich, N., Jaschek, R., Tanay, A. & Groner, Y. Dynamic combinatorial interactions of
702          RUNX1 and cooperating partners regulates megakaryocytic differentiation in cell line
703          models. *Blood* **117**, e1-14, doi:10.1182/blood-2010-07-295113 (2011).
704    29    Polak, R. *et al.* Autophagy inhibition as a potential future targeted therapy for ETV6-
705          RUNX1-driven B-cell precursor acute lymphoblastic leukemia. *Haematologica* **104**, 738-748,
706          doi:10.3324/haematol.2018.193631 (2019).
707    30    Wang, X. *et al.* Runx1 prevents wasting, myofibrillar disorganization, and autophagy of
708          skeletal muscle. *Genes Dev* **19**, 1715-1722, doi:10.1101/gad.1318305 (2005).
709    31    Tatetsu, H. *et al.* SALL4, the missing link between stem cells, development and cancer. *Gene*
710          **584**, 111-119, doi:10.1016/j.gene.2016.02.019 (2016).
711    32    Zhang, X., Yuan, X., Zhu, W., Qian, H. & Xu, W. SALL4: an emerging cancer biomarker and
712          target. *Cancer Lett* **357**, 55-62, doi:10.1016/j.canlet.2014.11.037 (2015).
713    33    Kong, N. R. *et al.* Zinc Finger Protein SALL4 Functions through an AT-Rich Motif to
714          Regulate Gene Expression. *Cell Rep* **34**, 108574, doi:10.1016/j.celrep.2020.108574 (2021).
715    34    Skene, P. J. & Henikoff, S. An efficient targeted nuclease strategy for high-resolution
716          mapping of DNA binding sites. *Elife* **6**, doi:10.7554/eLife.21856 (2017).
717    35    Baker, M. *1,500 scientists lift the lid on reproducibility*, <https://www.nature.com/news/1-
718          500-scientists-lift-the-lid-on-reproducibility-1.19970> (2016).
719    36    Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J. Transposition of
720          native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding
721          proteins and nucleosome position. *Nat Methods* **10**, 1213-1218, doi:10.1038/nmeth.2688
722          (2013).
723    37    Buenrostro, J. D., Wu, B., Chang, H. Y. & Greenleaf, W. J. ATAC-seq: A Method for
724          Assaying Chromatin Accessibility Genome-Wide. *Curr Protoc Mol Biol* **109**, 21 29 21-21 29
725          29, doi:10.1002/0471142727.mb2129s109 (2015).
726    38    Clark, S. J. *et al.* Genome-wide base-resolution mapping of DNA methylation in single cells
727          using single-cell bisulfite sequencing (scBS-seq). *Nat Protoc* **12**, 534-547,
728          doi:10.1038/nprot.2016.187 (2017).
729    39       (2015).
730    40    BBMap: A Fast, Accurate, Splice-Aware Aligner (USA, 2014).
731    41    Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
732          sequence data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170 (2014).
733    42    Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
734          *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).
735    43    Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-
736          2079, doi:10.1093/bioinformatics/btp352 (2009).
737    44    Ramirez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data
738          analysis. *Nucleic Acids Res* **44**, W160-165, doi:10.1093/nar/gkw257 (2016).
739    45    Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137,
740          doi:10.1186/gb-2008-9-9-r137 (2008).

741    46    Zang, C. *et al.* A clustering approach for identification of enriched domains from histone
742          modification ChIP-Seq data. *Bioinformatics* **25**, 1952-1958,
743          doi:10.1093/bioinformatics/btp340 (2009).
744    47    R: A Language and Environment for Statistical Computing (R Foundation for Statistical
745          Computing, Vienna, Austraia, 2021).
746    48    UpSetR: A More Scalable Alternative to Venn and Euler Diagrams for Visualizing
747          Intersecting Sets (2019).
748    49    VennDiagram: Generate High-Resolution Venn and Euler Plots (2018).
749    50    seqLogo: Sequence logos for DNA sequence alignments (2019).
750    51    ggplot2: Elegant Graphics for Data Analysis (2016).

751

# Figure 1

## a

**Peak-Motif Percentage for 10 ENCODE TF Datasets**

Legend: MAX, SPI1, CEBPB, ZBTB33, MAFF, RUNX1, ATF4, JUN, GATA1, MEIS2

X-axis: Peak-Motif %  (0.0, 0.2, 0.4, 0.6, 0.8, 1.0)
Y-axis: MACS2, Consensus

## b



Motif Position Bias from Weighted Peak Center for CEBPB
Motif Position Bias from Weighted Peak Center for JUN
Motif Position Bias from Weighted Peak Center for SPI1
Motif Position Bias from Weighted Peak Center for ZBTB33

Legend: Consensus, GEM, Genrich, HOMER, MACS2

## c
### CEBPB *De Novo* Motif Search - Consensus

HOMER

p-value $1*10^{-13461}$
1st ranked motif

p-value $1*10^{-472}$
2nd ranked motif

## d
### CEBPB *De Novo* Motif Search - MACS2

HOMER

p-value $1*10^{-27364}$
1st ranked motif

p-value $1*10^{-4486}$
2nd ranked motif

## e
### MAFF *De Novo* Motif Search - Consensus

HOMER

p-value $1*10^{-3639}$
1st ranked motif

p-value $1*10^{-886}$
2nd ranked motif

## f
### MAFF *De Novo* Motif Search - MACS2

HOMER

p-value $1*10^{-8602}$
1st ranked motif

p-value $1*10^{-1251}$
2nd ranked motif

# Figure 2



**a**

**b**

**c**

72,441    3832    6862

**Cut&Run**
(Peaks found in at least 2 replicates)

**ChIP-Seq**
(Union Peak Set)

**d** Directed SALL4 Motif Search

55.2% of union peaks have motif

**e** *De Novo* Motif Search - STREME

p-value 2.3*10$^{-15}$
2nd ranked motif

**f** *De Novo* Motif Search - HOMER

p-value 1*10$^{-160}$
3rd ranked motif

**g**

p = 5.0*10$^{-49}$

—Union Peak Set

Background Control Sequences

**Figure 3**

**a**

Input

Process

Output

**b**

Raw Fastq

Settings Table

*Folder #*   *01*   *02*   *08*   *22*

*11*
*12*
*13*
*14*

**c**



**d**



Data Input

Command Line Translation