

# Motif elucidation in ChIP-seq datasets with a knockout control

Danielle Denisko<sup>1,2,†</sup>, Coby Viner<sup>2,3,†</sup>, and Michael M. Hoffman<sup>1-4,\*</sup>

<sup>1</sup>Department of Medical Biophysics, University of Toronto, Toronto, ON, Canada

<sup>2</sup>Princess Margaret Cancer Centre, Toronto, ON, Canada

<sup>3</sup>Department of Computer Science, University of Toronto, Toronto, ON, Canada

<sup>4</sup>Vector Institute, Toronto, ON, Canada

October 21, 2019

## Abstract

Chromatin immunoprecipitation-sequencing (ChIP-seq) is widely used to find transcription factor binding sites, but suffers from various sources of noise. Knocking out the target factor mitigates noise by acting as a negative control. Paired wild-type and knockout experiments can generate improved motifs but require optimal differential analysis. We introduce *peaKO*—a method to automatically optimize motif analyses with knockout controls, which we compare to two other methods. *PeaKO* often improves elucidation of the target factor and highlights the benefits of knockout controls, which far outperform input controls. It is freely available at <https://peako.hoffmanlab.org>.

## Introduction

Transcription factors, often recognizing specific DNA motifs, control gene expression by binding to *cis*-regulatory DNA elements<sup>55</sup>. Accurate identification of transcription factor binding sites remains a challenge<sup>24</sup>, with experimental noise further compounding a difficult problem<sup>32</sup>. Improving motif models to better capture transcription factor binding affinities at each position of the binding site facilitates downstream analyses on gene-regulatory effects. Higher-quality motifs also promote the exclusion of spurious motifs, obviating costly experimental follow-up.

Chromatin immunoprecipitation-sequencing (ChIP-seq)<sup>29,60</sup> is a standard approach to locating DNA-binding protein and histone modification occupancy across the genome. Many steps of the ChIP-seq protocol can introduce noise, masking true biological signal and impeding downstream interpretation<sup>16,27,32,42,57</sup>. Poor antibody quality presents a major source of noise, characterized by low specificity to the target transcription factor or non-specific cross-reactivity. Cross-reactive antibodies often cause spurious pull-down of closely related transcription factor family members. Antibody clonality also contributes to antibody quality. Polyclonal antibodies tend to recognize multiple epitopes, which allows for more flexibility in binding to the desired transcription factor but at the cost of increasing background noise<sup>32</sup>.

To address issues of antibody quality, large consortia such as the Encyclopedia of DNA Elements (ENCODE) Project have established guidelines for validating antibodies through rigorous assessment of sensitivity and specificity<sup>22,42</sup>. Other considerable sources of technical noise include increased susceptibility to fragmentation in open chromatin regions<sup>4</sup>, and variations in sequencing efficiency of

\*Correspondence: [michael.hoffman@utoronto.ca](mailto:michael.hoffman@utoronto.ca)

<sup>†</sup>Danielle Denisko and Coby Viner contributed equally to this work

DNA segments arising from differences in base composition<sup>32</sup>. Downstream computational processing further reveals a different type of noise arising from contamination of peaks with zingers, motifs for non-targeted transcription factors<sup>72</sup>.

Additional control experiments can mitigate the effects of the aforementioned biases. Common types of controls include input and mock immunoprecipitation. Input control experiments isolate cross-linked and fragmented DNA without adding an antibody for pull down. Mock immunoprecipitation control experiments utilize a non-specific antibody, commonly immunoglobulin G (IgG)<sup>27,42</sup>, during the affinity purification step, instead of an antibody to the transcription factor. In theory, IgG mock experiments should better address technical noise since they more closely mimic the steps of the wild type (WT) ChIP protocol<sup>42</sup>. In practice, however, they suffer from a range of issues stemming from low yield of precipitated DNA<sup>32</sup>. Although the ENCODE Project<sup>22</sup> recommends the use of input controls, these experiments also suffer from limitations. Input can only capture biases in chromatin fragmentation and sequencing efficiencies, thus failing to capture the full extent of ChIP-seq technical noise.

Knockout (KO) control experiments present an attractive alternative to input and mock immunoprecipitation. In these experiments, mutations directed to the gene encoding the target transcription factor result in little to no expression of the transcription factor, prior to ChIP-seq. This preserves most steps of the ChIP protocol, including antibody affinity purification. Therefore, KO experiments can account for both antibody-related noise and biases in library preparation.

Common transcription factor KO constructs include CRISPR/Cas9-targeted mutations<sup>17</sup> and Cre/loxP conditional systems<sup>64,65</sup>. In downstream computational analyses, signal from the KO experiment serves as a negative set for subtraction from the WT positive set. Many pre-existing computational methods can use negative sets, typically input controls, to model background distributions<sup>58,68,75</sup>. For example, some peak calling tools, such as MACS2<sup>74</sup>, can perform discriminative peak calling. Most of these tools use the control set to set parameters of a background Poisson or negative binomial distribution<sup>5</sup> serving as a null for assessing the significance of WT peaks<sup>58</sup>.

Since KO controls better account for biases in WT data than input controls, optimizing methods for KO controls should improve the quality of results from downstream analyses. Indeed, as KO constructs become increasingly more accessible<sup>19</sup>, the need for optimal KO processing guidelines becomes more crucial. While some preliminary studies have investigated the use of KO controls<sup>34,48</sup>, further rigorous comparison of methods and establishment of a standard remain necessary.

To elucidate motifs when KO controls are available, we introduce a new method, *peaKO*. *PeaKO* combines two pipelines incorporating differential processing of WT and KO datasets at different stages. By comparing the rankings of a variety of known and *de novo* motifs, we highlight *peaKO*'s value for discovering and assessing binding motifs of WT/KO experiments, and *peaKO*'s applications in other differential contexts.

## Results

### PeaKO combines two differential analysis pipelines

Two steps of ChIP-seq computational processing allow for the subtraction of control signal from WT signal: peak calling and motif analysis. Therefore, we created two complementary pipelines, Pipeline A and Pipeline B, integrating the same software tools but selecting opposing steps to subtract matched KO signal from WT signal. (Figure 1A).

Pipeline A incorporates differential motif analysis through MEME-ChIP<sup>49,50</sup>. It focuses on the motif discovery algorithms MEME<sup>7,8</sup> and DREME<sup>6</sup>, and includes the motif enrichment algorithm CentriMo<sup>9,43</sup>. MEME-ChIP uses control peak sets for discriminative enrichment analysis<sup>49</sup>.

Instead of differential motif analyses, Pipeline B incorporates differential peak calling through MACS2<sup>74</sup>. MACS2 uses the control peak set to set the parameters of the background null distribution from which it calls significant peaks. Pipeline B drew inspiration from the knockout implemented normalization (KOIN) pipeline<sup>34</sup>.

Both pipelines conclude by executing CentriMo<sup>9,43</sup>. CentriMo's measure of motif central enrichment assesses the direct DNA binding of the enriched transcription factor<sup>9</sup>. Some aspects of CentriMo's output differ according to whether we choose differential<sup>43</sup> or non-differential<sup>9</sup> mode. Both pipelines, however, output a list of motifs ranked in order of increasing p-values. Ideally, the top motif should reflect the target factor in the underlying ChIP-seq experiment, although some circumstances may preclude this.

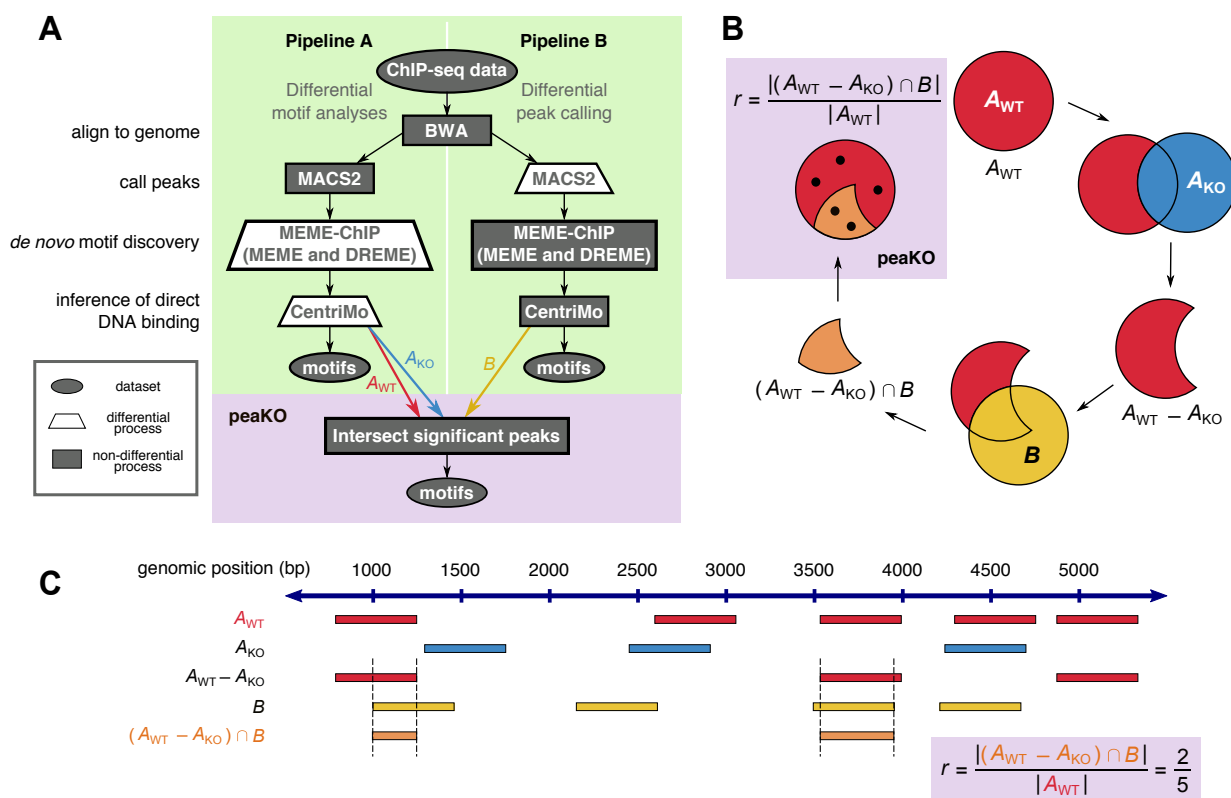
Each pipeline incorporates a unique approach to discriminative analysis. By modeling the peak background distribution using the negative control set, Pipeline B directly compares the position of read pileups between positive and negative datasets. In this model, we assume that read pileups shared between both datasets represent technical noise, while the remaining significant WT read pileups represent binding of the target transcription factor. Conversely, Pipeline A disregards the positional information of peaks and instead focuses on the position of the motif matches within the peaks. Pipeline A takes into account each peak's membership in the positive or negative set only when assessing the statistical significance of a motif. In Pipeline A, the simple motif discovery tool DREME compares the fraction of *de novo* motif matches in WT sequences to KO sequences. We assume that motifs more often located near peak centers in the WT dataset than in the KO dataset suggest associated binding events.

To select for motifs that both have consistent matches within peaks and fall within regions of significant read pileup, we combined both pipelines in a new way to develop *peaKO*. For each motif, *peaKO* computes the number of overlapping peaks between peak sets generated by both pipelines, with overlaps interpreted as genuine binding events (Figure 1B and Figure 1C; see Methods).

## PeaKO usually improves or maintains the best ranking of the known motif

To assess the performance of each method, we can first compare how well methods rank known canonical motifs of sequence-specific transcription factor datasets. We collected publicly available WT/KO paired ChIP-seq datasets for 8 sequence-specific transcription factors: ATF3<sup>76</sup>, ATF4<sup>26</sup>, CHOP<sup>26</sup>, GATA3<sup>70</sup>, MEF2D<sup>3</sup>, OCT4<sup>33</sup>, SRF<sup>66</sup>, and TEAD4<sup>30</sup> (Table 1). We evaluated our methods on these datasets, supplementing CentriMo with the collection of vertebrate motifs from the JASPAR 2016 database<sup>52</sup> (see Methods). Each transcription factor in our WT/KO datasets contains a corresponding motif within the JASPAR database. We used these JASPAR motifs as our gold-standard known motifs, and compared their rankings across methods. As a control, we processed the WT dataset alone through the same pipeline steps without any KO data.

In 5 out of 8 cases, *peaKO* improved or maintained the optimal rank relative to all other methods. *PeaKO* also always improved or maintained the rank relative to at least one other method (Figure 2). The total number of ranked motifs differed between experiments, which suggests *peaKO* may benefit analyses for a wide range of transcription factors with variable binding affinities. Of the other methods, Pipeline A performed the worst overall, as exemplified by rankings for the GATA3 (rank of 118) and ATF3 (rank of 240) datasets. Pipeline B performed similarly to the use of only WT data processed without controls, suggesting it benefits little from the control. *PeaKO* combines the best aspects of both types of differential analysis pipelines, limiting their deficiencies and highlighting their strengths. This generally leads to better rankings of known motifs.

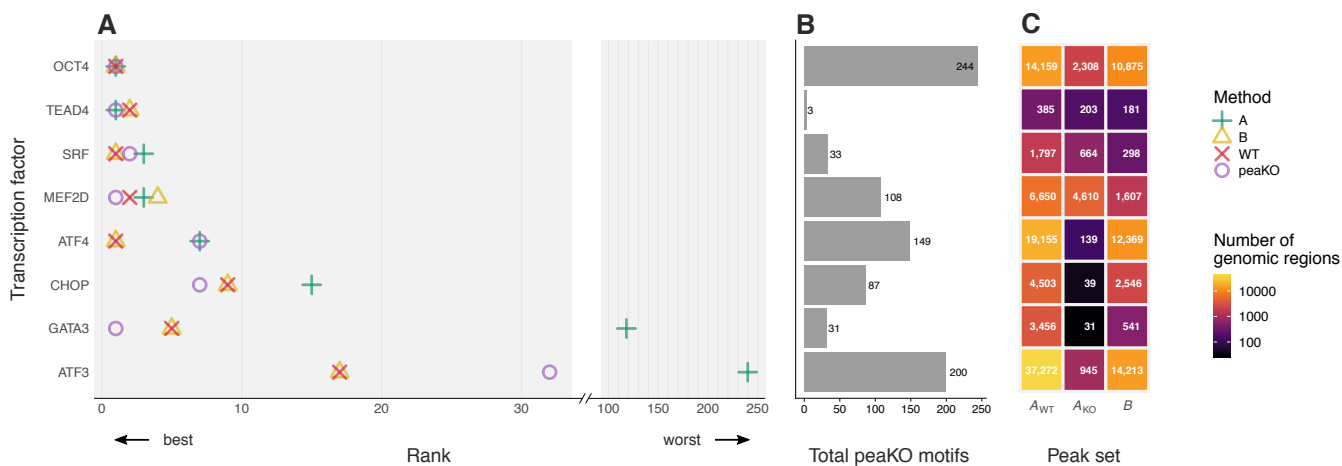


**Figure 1. Overview of Pipelines A and B, and peaKO.** (A) Pipelines A and B differ in their differential analysis steps. Each pipeline accepts both wild type (WT) and knockout (KO) ChIP-seq data as input. Pipeline A incorporates differential motif elucidation via MEME-ChIP<sup>50</sup>, whereas Pipeline B incorporates differential peak calling via MACS2<sup>74</sup>. Both pipelines produce a ranked list of motifs predicted as relevant to the ChIP-seq experiment by CentriMo<sup>9,43</sup>. PeaKO extracts significant peaks from CentriMo and computes a new score by which it ranks motifs. (B) PeaKO computes its ranking metric  $r$  through a series of set operations. PeaKO uses peak sets  $A_{WT}$  and  $A_{KO}$ , extracted from Pipeline A, and peak set  $B$ , extracted from Pipeline B. (C) A toy example illustrates the calculation of peaKO's score. Starting from the top row of peak set  $A_{WT}$  and moving downwards, we apply the peak set operations of  $r$  sequentially to identify regions satisfying the numerator criteria.

## De novo motifs consistently match known motifs

We investigated each method's ability to rank *de novo* motifs and assessed the similarity between *de novo* and known JASPAR motifs. For consistency, we pooled *de novo* motifs generated by each method (see Methods). We quantified similarity between *de novo* and known motifs using Tomtom<sup>25</sup>. We studied these methods on the same 8 WT/KO paired datasets used for our known motif analyses.

Usually, top *de novo* motifs more closely resembled the canonical motif across methods, resulting in most ranking near 1 (Figure 3). Conversely, motifs ranking lower tended to have fewer matches to the known motif, often not even matching the known motif at all. PeaKO generally followed this trend, but in a few exceptions, such as CHOP, OCT4, and ATF3, top motifs also sparsely matched the canonical motif. PeaKO might have found related, interacting factors, rather than the factor of interest. For example, many top *de novo* motifs reported by peaKO for the CHOP dataset closely matched the motif for ATF4, which interacts with CHOP<sup>26</sup>.



**Figure 2. Known transcription factor motifs elucidated by different methods.** Motifs originated from the JASPAR 2016 motif database<sup>52</sup>. Knockout datasets served as a control for differential analyses. **(A)** Each method ranked JASPAR database motifs based on their centrality within peak sets, as determined by CentriMo<sup>9,43</sup>. Ranks correspond to the ChIPped transcription factor’s known motif (Table 2). **(B)** Total number of motifs assessed by peakKO. **(C)** The number of peaks found by each method varies across peak sets.

## PeaKO teases apart similar GATA family motifs

We delved deeper into our GATA3 results, for which peaKO outperformed all other methods. GATA3 belongs to the family of GATA factors, all of which bind GATA-containing sequences<sup>54</sup>. Despite having similar motifs, each GATA factor plays a distinct role and usually does not interact with the others<sup>69</sup>.

Distinguishing the targeted motif among GATA factors and other large transcription factor families often presents a challenge. Minor differences in position weight matrices (PWMs)<sup>14</sup> can cause major differences in genome-wide transcription factor binding sites<sup>39</sup>. Understanding the downstream effects of transcription factor binding necessitates pulling apart these intricacies in motif preferences.

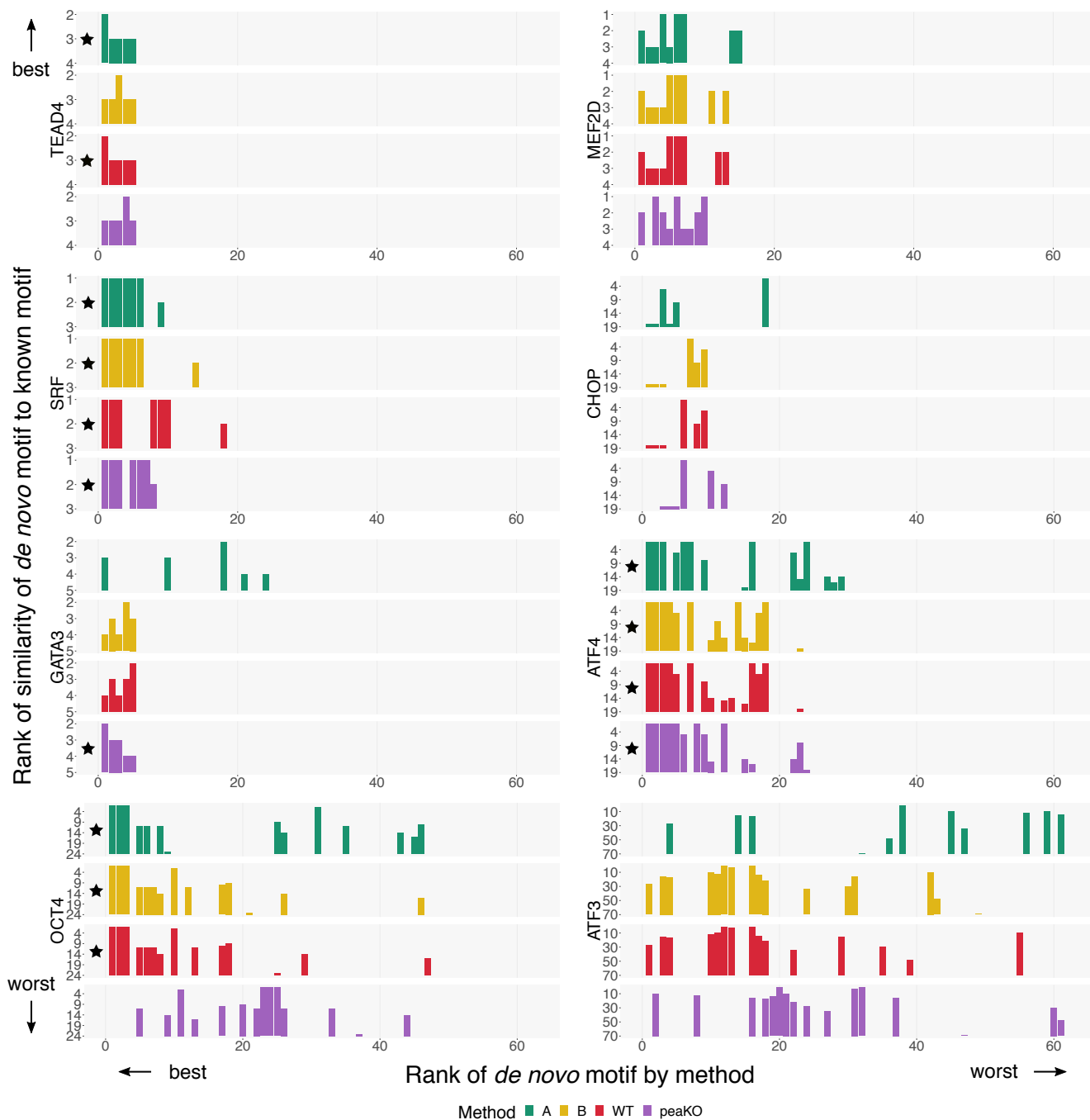
CentriMo results across both pipelines further reinforced the difficulty of distinguishing these motifs (Figure 4). Pipeline B identified closely related GATA family members with ranks 1–4, above the desired fifth-ranked GATA3 motif. Pipeline A proved less promising, failing to rank any GATA members within its top 10 motifs. Furthermore, none of the top Pipeline A motifs appeared centrally enriched within WT peaks. Instead, we observed a uniform distribution among the WT peak set and a series of stochastic, sharp peaks among the KO peak set, likely representing inflated probabilities due to low sample size.

Despite the difficulties affecting Pipeline B, peaKO draws on its ability to detect GATA family members, and surpasses it by ranking GATA3 first. Thus, we find peaKO achieves specificity in ranking motifs in the presence of many similar motifs.

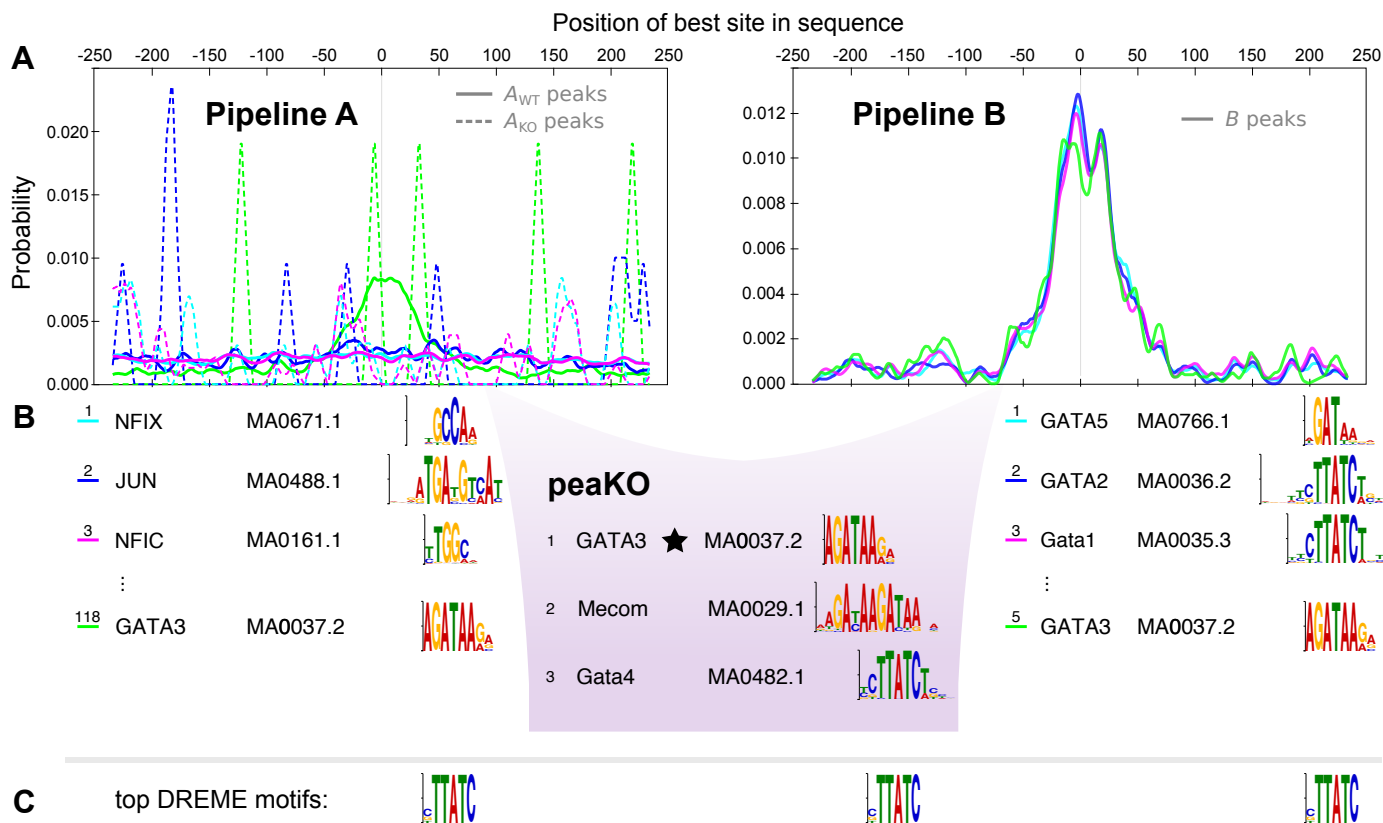
## Low-quality datasets account for poor rankings across methods

In a few cases, peaKO performed worse than the other methods at ranking the canonical motif (Figure 2). In particular, we observed a large spread in rankings across methods for ATF3 (ranging from rank 17 to rank 240). We found central enrichment of the canonical ATF3 motif in the KO peak set, as depicted by Pipeline A’s CentriMo results (Figure 5). This central enrichment appears even more prominent than that in the WT peak set.

Although CentriMo probabilities depend on the total number of peaks in each set, and a relatively low number of peaks in the control set can inflate these probabilities, we expect non-specific matches



**Figure 3. Similarity of discovered *de novo* motifs to canonical JASPAR motifs across 4 methods.** For 8 transcription factors (Table 2), we ran 4 methods (green: Pipeline A, yellow: Pipeline B, red: WT alone, purple: peaKO) on a pooled set of *de novo* motifs generated by MEME<sup>7,8</sup> and DREME<sup>6</sup>. Each method generated a ranking of *de novo* motifs. For each of these motifs, we quantified similarity to the known motif using Tomtom<sup>25</sup>. To emphasize strong matches to known motifs, the provided ranks lie in descending order, with the best (rank 1) motif, at the top. In some cases, the best rank achieved by the match does not reach 1, as reflected by a greater lower limit. Black stars: methods achieving the best possible rank across both ranking schemes within each experiment.



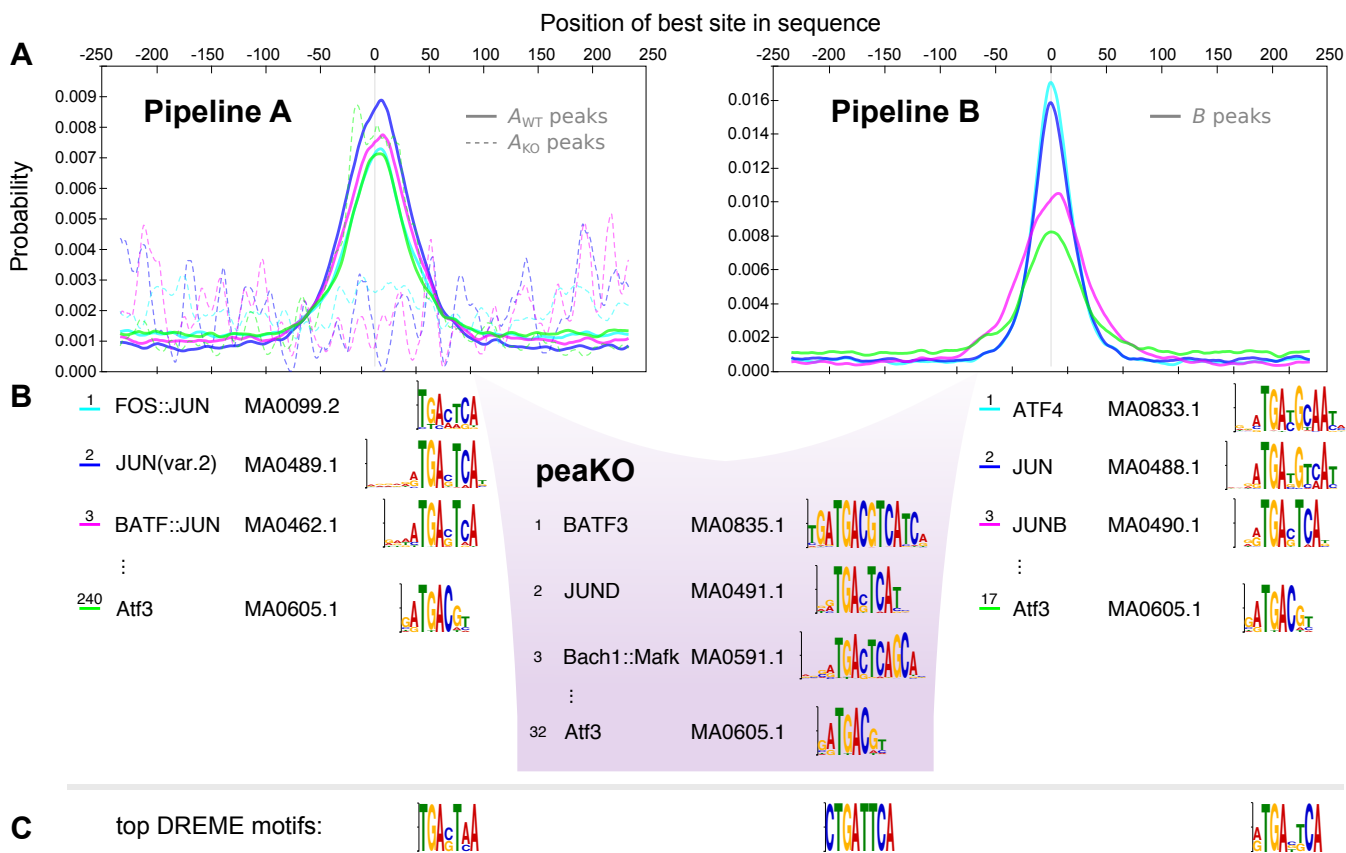
**Figure 4. PeaKO ranks the GATA3 motif above other GATA transcription factor family motifs.** (A) CentriMo<sup>9,43</sup> probability plots depict enrichment of the top 3 motifs from each method, along with the GATA3 motif, within peak sets. (B) Motifs resulting from each pipeline and PeaKO lie beneath associated CentriMo plots. Motifs and corresponding sequence logos<sup>63</sup> originate from JASPAR 2016<sup>52</sup>. Capitalization is as it occurs in JASPAR. The information content of bases in the sequence logos ranges from 0 bits to 2 bits. The black star denotes achieving the best rank of the GATA3 motif. (C) Top DREME<sup>6</sup> motifs with length greater than 5 bp, for comparison. In this case, all three motifs are identical.

to generate a uniform background distribution rather than a distinctive centrally-enriched pattern<sup>9,43</sup>. Accordingly, ATF3 enrichment deviates substantially from our expectations and suggests issues with the underlying KO ChIP experiment. This likely explains the poor rankings of ATF3 across methods, including PeaKO.

## Knockout-controlled analyses consistently improve motif elucidation

To investigate whether KO controls would better approximate WT ChIP-seq experimental noise than input controls, we used input controls to repeat our analyses. We ran our methods on MEF2D, OCT4, and TEAD4 datasets, which contained input controls (Table 1), by applying the same procedures but using only the input dataset for differential analysis steps.

Using an input control instead of a KO control usually worsened the ranking of the known motif, as observed by an overall shift across methods toward poorer rankings (Figure 6A). In *de novo* motif analyses with input controls, top-ranked motifs tended to have slightly poorer matches to known motifs across methods, as compared with KO controls (Figure 6B). As in WT/KO analyses of OCT4,

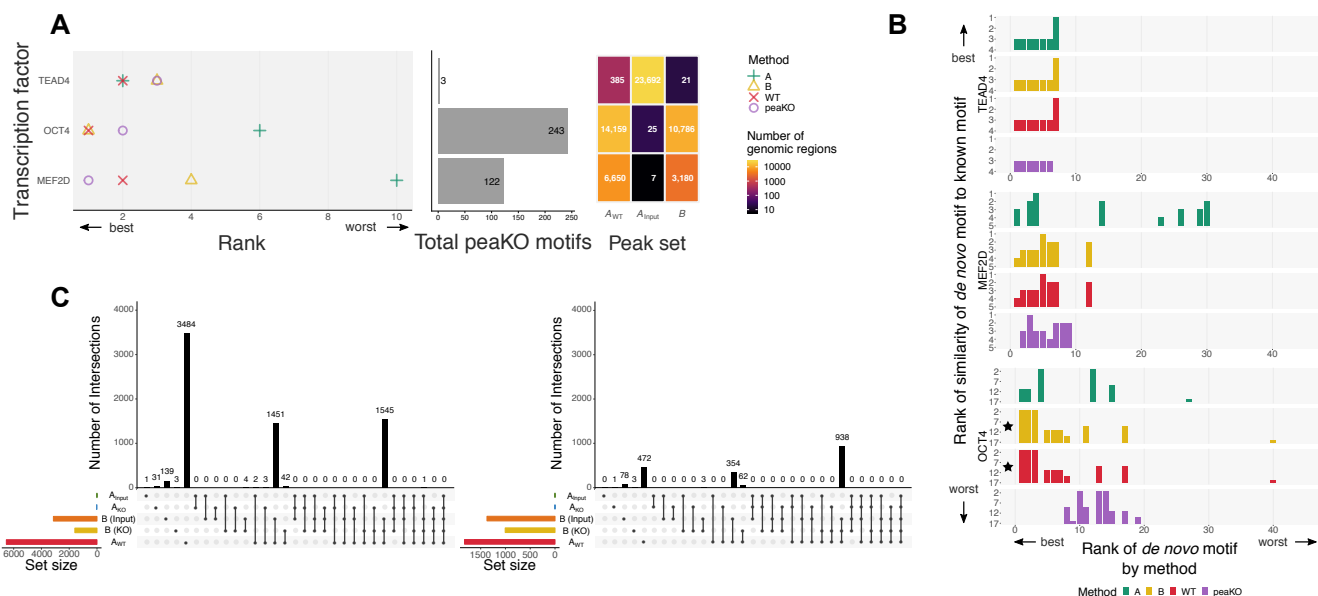


**Figure 5. The ATF3 motif is centrally enriched in the ATF3 knockout dataset.** (A) CentriMo<sup>9,43</sup> probability plots depict enrichment of the top 3 motifs from each method, along with the ATF3 motif, within peak sets. (B) Motifs resulting from each pipeline and peaKO lie beneath associated CentriMo plots. Motifs and corresponding sequence logos<sup>63</sup> originate from JASPAR 2016<sup>52</sup>. Capitalization is as it occurs in JASPAR. Information content of bases underlying motifs range from 0 bits to 2 bits. (C) Top DREME<sup>6</sup> motifs with length greater than 5 bp, for comparison.

we observed sparsity in top-ranked peaKO motifs matching the known motif. This could point to low affinity of the antibody to the target factor or other types of noise affecting primarily the WT set. Indeed, input experiments yielded even fewer significant peaks from CentriMo than KO experiments (Figure 6C).

Overall, using input controls instead of KO controls led to poorer rankings across methods. Although peaKO did not outperform the other methods using only input, it generally performed similarly, suggesting utility in other differential applications.





**Figure 6. Ranks of known and discovered motifs using input controls.** (A) Ranks of known JASPAR<sup>52</sup> motifs across methods for each ChIP-seq experiment (Table 2). Input datasets served as a control in differential analysis steps. (B) We plotted ranks of *de novo* motifs discovered by MEME<sup>7,8</sup> and DREME<sup>6</sup> against their similarity to the known JASPAR motif, as quantified by Tomtom<sup>25</sup>. We compared queried motifs against the JASPAR 2016 target motif database. Black stars: methods achieving the best possible rank across both ranking schemes within each experiment. (C) UpSet plot<sup>44</sup> of overlap between MEF2D peak sets generated by Intervene<sup>31</sup> (left) for all motifs and (right) for the MEF2D motif only. For Pipeline B peak sets, parentheses indicate the type of negative control used for peak calling: input or knockout (KO).

## Discussion

Increased accessibility of KO experiments presents a need for standardized computational processing workflows. With KO data, peaKO's dual pipeline approach generally outperformed each pipeline alone when ranking the known motifs. This holds true even in challenging cases, such as distinguishing among large transcription factor families with shared core motifs. Applying our methods to datasets containing both input and KO controls demonstrates the superiority of KO controls for motif elucidation.

We observed a common theme throughout our analyses pertaining to the characteristic performance of each pipeline alone. When tasked with ranking the known motif, Pipeline A generally produced inferior rankings, especially for ATF3 and GATA3 (ranks > 100) and, to a lesser extent, CHOP (rank 15). We could only attribute this to poor experimental quality for ATF3. The significance of differential mode CentriMo p-values, calculated using Fisher's exact test, appears closely linked to the relative size of each peak set. Both CHOP and GATA3 KO control sets had fewer than 50 KO peaks (Figure 2), which might account for Pipeline A's poor performance.

Pipeline B suffered from a different issue: it ranked known motifs almost identically to WT processing alone, without any controls. Since the sole difference between Pipeline B and WT-only processing lies in the peak calling step, identical rankings indicate the sufficiency of constructing the background distribution with WT-derived values alone. Differential peak calling with KO controls does, however, reduce the size of the WT peak set. Perhaps this improves an already specific peak set such that the improvement is undetectable when ranking known motifs. Nonetheless, rankings differ in some cases and *de novo* motif analyses reveal differences between Pipeline B and WT-only processing. Overall,

both pipelines show strengths in specific contexts, which *peaKO* emphasizes.

Some of our methods ranked the motif of interest less favorably than other GATA family member motifs. Finding the general familial motif could prove sufficient in some cases<sup>61</sup>. Nonetheless, finding the specific motif helps with understanding the roles of individual transcription factors. GATA family members share a common core motif, yet each have distinct and detectable binding preferences that contribute to their diversity in genome-wide occupancy and function<sup>54</sup>.

For OCT4 (also known as POU5F1), we selected the Pou5f1::Sox2 motif ([MA0142.1](#)). SOX2, like OCT4, regulates pluripotency in embryonic stem cells<sup>73</sup>. The two transcription factors often act together to regulate gene expression by forming a complex and co-binding to DNA<sup>1</sup>. Here, however, the heterodimer motif differs substantially from the OCT4 motif alone, as it additionally contains a SOX2 motif<sup>1</sup>. We chose to use the heterodimer motif in assessing our methods because the authors of the study that generated the OCT4 dataset found a substantially larger proportion of peaks containing the heterodimer (44.0%) as compared to the monomer (20.6%)<sup>33</sup>. Upon re-running our analyses using the monomer motif instead, we found poorer rankings across methods, as expected from this imbalance of motif types in peaks (see <https://doi.org/10.5281/zenodo.3338330>). Higher occupancy of the heterodimer form, however, does not preclude the transcription factor from binding DNA in its monomer form. Although all methods found the heterodimer motif as the top rank, deciding upon which motif form to use and how it affects downstream processing would benefit from further exploration.

Our use of cross-species PWMs potentially limits our findings. We used motifs from the JASPAR vertebrate collection interchangeably where the known motif did not always originate from the same species as our ChIP-seq datasets (see Tables 1 and 2). Recently, Lambert et al<sup>41</sup> found that, contrary to commonly held belief, extensive motif diversification among orthologous transcription factors occurs quickly as species diverge. Additionally, PWMs<sup>14</sup> themselves, while providing the most commonly used motif model<sup>12,36</sup>, may not sufficiently capture nuanced binding differences<sup>20,36</sup>.

Lastly, we used *peaKO* along with our other methods to assess the benefit of KO controls over input, suggesting that *peaKO* may prove useful for other non-WT/KO differential contexts. CRISPR epitope tagging ChIP-seq (CETCh-seq), which involves the insertion and expression of FLAG epitope tags on the target transcription factor<sup>62</sup>, presents one alternative differential context which may gain from *peaKO*. CETCh-seq provides a substantial advantage over traditional ChIP-seq because it only requires one high-quality monoclonal antibody recognizing the FLAG antigen across any number of transcription factor experiments. Preliminary analyses using CETCh-seq datasets revealed challenges arising from unexpected signal from a shared control of ChIP-seq in an untagged cell line. Further work should investigate the role of CETCh-seq controls and how they integrate with *peaKO*.

Similar considerations for the proper use of control sets could also apply to combining replicates. Combining negative control replicates with the irreproducible discovery rate (IDR) framework<sup>47</sup> may pose problems considering that these datasets represent noise rather than a full range across true signal and noise. This may present an issue as IDR's underlying copula mixture model assumes the existence of an inflection point within the dataset marking the transition between true signal and noise<sup>47</sup>.

## Conclusion

We present *peaKO*, a free and publicly available tool for ChIP-seq motif analyses with KO controls (<https://peako.hoffmanlab.org>). *PeaKO* improves over two kinds of differential processing in ranking the motif of interest. We anticipate that *peaKO* will prove useful in identifying motifs of novel transcription factors with available KO controls. We hope this will encourage both greater collection and wider usage of knockout datasets.

## Methods

### Overview of ChIP-seq processing and analysis methods

ChIP-seq processing follows this overarching path:

1. subject sequenced reads to trimming and quality control assessment;
2. align reads to a reference genome;
3. call peaks according to significant read pileups; and
4. elucidate *de novo* motifs and assess peaks for evidence of direct DNA binding.

For some methods, steps 3 and 4 can incorporate information from control datasets. We constructed two pipelines to compare differential analyses in both of these steps (Figure 1A).

In Pipeline A, we perform differential analysis with MEME-ChIP<sup>49,50</sup>. MEME-ChIP uses the *de novo* motif elucidation tools MEME<sup>7,8</sup> and DREME<sup>6</sup>, and assesses the central enrichment of motifs in peaks via CentriMo<sup>9,43</sup>. CentriMo ranks motifs according to multiple-testing corrected binomial p-values (non-differential mode)<sup>9</sup> or Fisher's exact test p-values (differential mode)<sup>43</sup>.

In Pipeline B, we perform differential peak calling through MACS2<sup>74</sup>. While Pipeline B draws inspiration from the KOIN pipeline<sup>34</sup>, it does not incorporate the HOMER `makeTagDirectory` or `annotatePeaks`<sup>28</sup> steps. We replaced HOMER motif tools<sup>28</sup> with those from the MEME Suite<sup>10,11</sup>. Both Pipelines A and B incorporate identical pre-processing and alignment steps, described later. Since both pipelines employ CentriMo in their last step, they generate a list of ranked motifs with predicted association to the ChIP-seq experiment.

### PeaKO: motivation and score

Differential peak calling and differential motif analysis address the same problem of noise removal, albeit in distinct ways. Therefore, we surmised that by combining the two approaches, the results from each pipeline could complement and strengthen one another. CentriMo produces a ranked list of motifs, and each motif has an associated peak set containing a centered window enriched for that motif. We reasoned that motifs with a large proportion of peaks shared between both pipelines are likely relevant to the ChIP-seq experiment. We then created a metric that captures this.

PeaKO takes as input the CentriMo output of each pipeline. We modified CentriMo code to output negative control set peaks associated with each motif in differential mode, since current versions only output positive peaks. These changes are merged into the CentriMo source repository and the MEME Suite's next major release will include them. From the CentriMo results `peaKO` filters out motifs with multiple-testing corrected p-values  $> 0.1$ .

PeaKO computes a ranking metric  $r$  that represents the proportion of high-quality  $A_{WT}$  peaks found in set  $B$  but not in set  $A_{KO}$ . To do this, `peaKO` calculates the overlap between peak sets  $A_{WT}$  and  $A_{KO}$  from Pipeline A, and peak set  $B$  from Pipeline B through a series of set operations:

$$r = \frac{|(A_{WT} - A_{KO}) \cap B|}{|A_{WT}|}.$$

PeaKO implements these operations using `pybedtools` (version 0.7.7; BEDTools version 2.26.0)<sup>18,59</sup>. First, `peaKO` removes any  $A_{WT}$  peak overlapping at least 1 bp of a  $A_{KO}$  peak (`pybedtools subtract -A`; Figure 1B). Second, `peaKO` finds regions overlapping by at least 1 bp between remaining  $A_{WT}$  peaks and  $B$  peaks (`pybedtools intersect -wa`). Third, `peaKO`

applies `pybedtools merge` with default settings to overlapping regions, which merges identical regions and ensures that the ranking metric  $r$  has a maximum value of 1. PeaKO's final output consists of a list of motifs ranked according to this metric.

## Datasets

We analyzed a total of 8 publicly available ChIP-seq experiment datasets with KO controls (Table 1). We selected two datasets (GATA3 and SRF) from Krebs et al<sup>34</sup>, while we selected the remainder by searching for KO-associated ChIP-seq datasets on Gene Expression Omnibus (GEO)<sup>21</sup>. We accessed datasets through GEO, except for the SRF dataset<sup>28</sup>, available on Zenodo (<https://doi.org/10.5281/zenodo.3405482>). ATF3 experiments come from human tissue, while the other experiments come from mouse tissue.

**Table 1. ChIP-seq datasets used, with associated GEO accession numbers (where applicable) and number of replicates.**

Factor	GEO	Reference	Wild type	Knockout	Input
ATF3	<a href="#">GSE74355</a>	Zhao et al <sup>76</sup>	1	1	0
ATF4	<a href="#">GSE35681</a>	Han et al <sup>26</sup>	3	3	0
CHOP	<a href="#">GSE35681</a>	Han et al <sup>26</sup>	1	1	0
GATA3	<a href="#">GSE20898</a>	Wei et al <sup>70</sup>	1	1	0
MEF2D	<a href="#">GSE61391</a>	Andzelm et al <sup>3</sup>	3	1	3
OCT4	<a href="#">GSE87822</a>	King and Klose <sup>33</sup>	3	3	1
SRF	—	Sullivan et al <sup>66</sup>	1	1	0 <sup>a</sup>
TEAD4	<a href="#">GSE82190</a>	Joshi et al <sup>30</sup>	1	1	1

<sup>a</sup> We excluded the available dataset because it came from a different, older DNA sequencer and lacked quality scores.

## Motifs

We downloaded the collection of vertebrate motifs in MEME format<sup>11</sup> from the JASPAR CORE 2016 motif database, which consists of curated PWMs derived from *in vivo* and *in vitro* methods<sup>52</sup>.

We defined each canonical motif from the JASPAR collection as the motif matching the target transcription factor except in two cases: OCT4 and CHOP (Table 2). In both cases, we instead chose motifs derived from their common heterodimer complex forms. CHOP or DDIT3 likely binds DNA as an obligate multimer<sup>40,56</sup>, so we used `Ddit3::Cebpa` ([MA0019.1](#)). The CHOP monomer motif closely resembles its C/EBP $\alpha$  heterodimer motif, relative to its Cis-BP (version 1.02)<sup>71</sup> DDIT3 motif (`T025314_1.02`, derived from `HOCOMOCO`<sup>37</sup>). For OCT4, we used the `Pou5f1::Sox2` motif ([MA0142.1](#); see Discussion).

We provided motifs to CentriMo<sup>9,43</sup> for central enrichment analyses and to Tomtom<sup>25</sup> for similarity assessments.

## Pre-processing, alignment, and peak calling

Before alignment, we trimmed adapter sequences with TrimGalore! (version 0.4.1)<sup>35</sup> which uses Cutadapt (version 1.8.3)<sup>51</sup>. We performed quality control using FastQC (version 0.11.5)<sup>2</sup>. We used Picard's `FixMateInformation` and `AddOrReplaceReadsGroups` (version 2.10.5)<sup>15</sup> and GATK's `PrintReads` (version 3.6)<sup>53</sup> to prevent GATK errors. We then aligned reads to

**Table 2. Known JASPAR CORE 2016<sup>52</sup> vertebrate motifs.** Sentence case motif names designate mouse transcription factor motifs, while full upper case names designate human motifs. Double colons designate heterodimer motifs.

Common name	JASPAR ID	Motif name
ATF3	MA0605.1	Atf3
ATF4	MA0833.1	ATF4
CHOP	MA0019.1	Ddit3::Cebpa
GATA3	MA0037.2	GATA3
MEF2D	MA0773.1	MEF2D
OCT4	MA0142.1	Pou5f1::Sox2
SRF	MA0083.3	SRF
TEAD4	MA0809.1	TEAD4

GRCm38/mm10 or GRCh38/hg38 with BWA bwa-a1n (version 0.7.15)<sup>45</sup> (as recommended<sup>46</sup>, since some datasets have reads  $\ll$  70 bp), using Sambamba (version 0.6.6)<sup>67</sup> for post-processing.

Next, we called peaks using MACS2 (version 2.0.10)<sup>74</sup> with parameters `-q 0.05`. In Pipeline A, we called WT and KO peaks separately. In Pipeline B, we provided the KO dataset as a control to the WT dataset during peak calling (parameter `-c`), resulting in a single set of peaks.

## Combining replicates

For MEF2D, OCT4, and TEAD4 experiments which consist of biological replicates (see Table 1), we processed replicates using the **ENCODE Transcription Factor and Histone ChIP-seq processing pipeline**<sup>38</sup>. The ENCODE pipeline replaces the pre-processing, alignment, and peak calling steps described earlier. We chose default parameters for punctate (narrow peak) binding experiments in all steps. Instead of a q-value threshold, this pipeline caps the number of peaks ( $n = 500\,000$ ) to ensure that the IDR framework<sup>47</sup> can analyze a sufficient number of peaks across a full spectrum. IDR combines peaks across replicates based on the assumption that strong peaks shared across replicates represent true binding events, while weak, one-off peaks represent noise. To emulate the first steps of Pipeline A and Pipeline B, we either ran the ENCODE pipeline on WT replicates and KO replicates separately (for Pipeline A), or we ran the ENCODE pipeline on all WT and KO replicates simultaneously, setting KO replicates as controls (for Pipeline B). For downstream motif analyses, we used the combined “optimal” peak sets output by IDR.

## Motif analyses with MEME-ChIP

In both pipelines, we employed MEME-ChIP<sup>49,50</sup> from the MEME Suite<sup>10,11</sup> for motif analysis. We used MEME-ChIP version 4.12.0, except for CentriMo, which we compiled from version 4.11.2 and modified to output negative sequences. MEME-ChIP performs motif discovery with complementary algorithms MEME<sup>7,8</sup> and DREME<sup>6</sup>, and motif enrichment with CentriMo<sup>9,43</sup>.

We extended MACS2 narrowPeak regions equidistantly from peak summits to create a uniform set of 500 bp centered peaks<sup>49</sup>. We then extracted underlying genomic sequences using BEDTools s1op (version 2.23.0)<sup>59</sup> from a repeat-masked genome. We masked the genome with Tandem Repeats Finder (TRF) (version 4.09)<sup>13</sup> with options `-h -m -ngs` and parameters 2 7 7 80 10 50 500 for mouse (as done originally by Benson<sup>13</sup>), and options 2 5 5 80 10 30 200 for human (as recommended by Frith et al<sup>23</sup>).

In Pipeline A, we provided the negative control set in addition to the WT set, running MEME, DREME, and CentriMo in differential mode. In ranking known motifs, we ran CentriMo providing only JASPAR database motifs. In ranking *de novo* motifs, we ran CentriMo providing only MEME and DREME motifs.

## Pooling *de novo* motifs

Each run of MEME or DREME creates new and globally non-unique identifiers for output motifs. This leads to recurring identifiers that refer to different motifs across multiple runs. To consolidate identifiers across multiple MEME and DREME runs, we modified identifiers to reflect the pipeline from which they originate. We then pooled *de novo* motifs across methods and re-ran the CentriMo step of each pipeline, providing the pooled database, allowing for accurate comparisons.

## Assessing similarity of *de novo* motifs to known motifs

For each experiment, we quantified the similarity of *de novo* motifs to the known JASPAR motif using Tomtom<sup>25</sup>. Tomtom compared the *de novo* motifs to the JASPAR motif database through ungapped alignment across columns<sup>25</sup>. Tomtom generated a list of known motif matches, ranked by increasing Bonferroni-corrected p-values. An exact match between a *de novo* motif and a JASPAR motif would result in the JASPAR motif's ranking first in this list of matches.

## Comparing input to knockout controls

For experiments with associated input controls, we re-ran our known motif and *de novo* motif analyses swapping out KO datasets for input datasets. We compared peaks between sets using UpSet (version 1.4.0) plots<sup>44</sup>, via Intervene (version 0.6.2)<sup>31</sup>, which calculates genomic region overlaps with BEDTools (version 2.26.0)<sup>59</sup>.

## Availability of data and materials

PeaKO is available at <https://peako.hoffmanlab.org> with Python source code for peaKO and both pipelines at: <https://github.com/hoffmangroup/peako>. Persistent availability is ensured by Zenodo, in which we have deposited the current version of our code (<https://doi.org/10.5281/zenodo.3338324>), its downstream CentriMo and peaKO outputs (<https://doi.org/10.5281/zenodo.3338330>), and our changes to the CentriMo source code and the Linux x86-64 binary that we used (<https://doi.org/10.5281/zenodo.3356995>). All source code is licensed under a [GNU General Public License, version 3 \(GPLv3\)](https://www.gnu.org/licenses/gpl-3.0.html), except for CentriMo, which retains its original license.

## Competing interests

The authors declare that they have no competing interests.

## Authors' contributions

Conceptualization, C.V.; Data Curation, D.D. and C.V.; Methodology, D.D., C.V., and M.M.H.; Software, D.D. and C.V.; Visualization, D.D., C.V., and M.M.H.; Writing — Original Draft, D.D.; Writing — Review & Editing, D.D., C.V., and M.M.H.; Resources, M.M.H.; Funding Acquisition, M.M.H.; Supervision, C.V. and M.M.H.

## Acknowledgments

We thank Christopher K. Glass (ORCID: 0000-0003-4344-3592) and Verena Link (ORCID: 0000-0002-3207-312X) for providing SRF datasets. We thank Carl Virtanen (ORCID: 0000-0002-2174-846X), Zhibin Lu (ORCID: 0000-0001-6281-1413), and Qun Jin (Bioinformatics and High Performance Computing Core, University Health Network).

## Funding

This work was supported by the Natural Sciences and Engineering Research Council of Canada (Alexander Graham Bell Canada Graduate Scholarships to D.D. and C.V.), the Canadian Institutes of Health Research (201512MSH-360970 to M.M.H. and Undergraduate Summer Studentship Award to D.D.), the Ontario Ministry of Training, Colleges and Universities (Ontario Graduate Scholarships to D.D. and C.V.), the Ontario Ministry of Research, Innovation and Science (ER-15-11-223 to M.M.H.), the University of Toronto Undergraduate Research Opportunities Program (to D.D.), and the Princess Margaret Cancer Foundation.

## References

- [1] Aksoy I, Jauch R, Chen J, Dyla M, Divakar U, Bogu GK, Teo R, Leng Ng CK, Herath W, Lili S, Hutchins AP, Robson P, Kolatkar PR, Stanton LW (2013) Oct4 switches partnering from Sox2 to Sox17 to reinterpret the enhancer code and specify endoderm. *EMBO J* 32(7):938–953, <https://doi.org/10.1038/emboj.2013.31>
- [2] Andrews S (2018) FastQC: a quality control tool for high throughput sequence data. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
- [3] Andzelm MM, Cherry TJ, Harmin DA, Boeke AC, Lee C, Hemberg M, Pawlyk B, Malik AN, Flavell SW, Sandberg MA, Raviola E, Greenberg ME (2015) MEF2D drives photoreceptor development through a genome-wide competition for tissue-specific enhancers. *Neuron* 86(1):247–263, <https://doi.org/10.1016/j.neuron.2015.02.038>
- [4] Auerbach RK, Euskirchen G, Rozowsky J, Lamarre-Vincent N, Moqtaderi Z, Lefrançois P, Struhl K, Gerstein M, Snyder M (2009) Mapping accessible chromatin regions using Sono-Seq. *Proc Natl Acad Sci USA* 106(35):14926–14931, <https://doi.org/10.1073/pnas.0905443106>
- [5] Bailey T, Krajewski P, Ladunga I, Lefebvre C, Li Q, Liu T, Madrigal P, Taslim C, Zhang J (2013) Practical guidelines for the comprehensive analysis of ChIP-seq data. *PLOS Comput Biol* 9(11):e1003326, <https://doi.org/10.1371/journal.pcbi.1003326>
- [6] Bailey TL (2011) DREME: motif discovery in transcription factor ChIP-seq data. *Bioinformatics* 27(12):1653–1659, <https://doi.org/10.1093/bioinformatics/btr261>
- [7] Bailey TL, Elkan C (1994) Fitting a mixture model by expectation maximization to discover motifs in biopolymers. In: Altman R, Brutlag D, Karp P, Lathrop R, Searls D (eds) *Proc Int Conf Intell Syst Mol Biol*, vol 2, pp 28–36
- [8] Bailey TL, Elkan C (1995) Unsupervised learning of multiple motifs in biopolymers using expectation maximization. *Mach Learn* 21(1–2):51–80, <https://doi.org/10.1007/BF00993379>
- [9] Bailey TL, Machanick P (2012) Inferring direct DNA binding from ChIP-seq. *Nucleic Acids Res* 40(17):e128, <https://doi.org/10.1093/nar/gks433>
- [10] Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW, Noble WS (2009) MEME Suite: tools for motif discovery and searching. *Nucleic Acids Res* 37(Web Server Issue):W202–W208, <https://doi.org/10.1093/nar/gkp335>

- [11] Bailey TL, Johnson J, Grant CE, Noble WS (2015) The MEME Suite. *Nucleic Acids Res* 43(W1):W39–W49, <https://doi.org/10.1093/nar/gkv416>
- [12] Benos PV, Bulyk ML, Stormo GD (2002) Additivity in protein-DNA interactions: how good an approximation is it? *Nucleic Acids Res* 30(20):4442–4451, <https://doi.org/10.1093/nar/gkf578>
- [13] Benson G (1999) Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res* 27(2):573–580, <https://doi.org/10.1093/nar/27.2.573>
- [14] Berg OG, von Hippel PH (1987) Selection of DNA binding sites by regulatory proteins. *J Mol Biol* 193(4):723–743, [https://doi.org/10.1016/0022-2836\(87\)90354-8](https://doi.org/10.1016/0022-2836(87)90354-8)
- [15] Broad Institute (2015) Picard. <http://broadinstitute.github.io/picard>
- [16] Chen Y, Negre N, Li Q, Mieczkowska JO, Slattery M, Liu T, Zhang Y, Kim TK, He HH, Zieba J, Ruan Y, Bickel PJ, Myers RM, Wold BJ, White KP, Lieb JD, Liu XS (2012) Systematic evaluation of factors influencing ChIP-seq fidelity. *Nat Methods* 9(6):609–614, <https://doi.org/10.1038/nmeth.1985>
- [17] Cong L, Ran FA, Cox D, Lin S, Barretto R, Habib N, Hsu PD, Wu X, Jiang W, Marraffini LA, Zhang F (2013) Multiplex genome engineering using CRISPR/Cas systems. *Science* 339(6121):819–823, <https://doi.org/10.1126/science.1231143>
- [18] Dale RK, Pedersen BS, Quinlan AR (2011) Pybedtools: a flexible Python library for manipulating genomic datasets and annotations. *Bioinformatics* 27(24):3423–3424, <https://doi.org/10.1093/bioinformatics/btr539>
- [19] Doudna JA, Charpentier E (2014) The new frontier of genome engineering with CRISPR-Cas9. *Science* 346(6213):1258096, <https://doi.org/10.1126/science.1258096>
- [20] Dror I, Rohs R, Mandel-Gutfreund Y (2016) How motif environment influences transcription factor search dynamics: finding a needle in a haystack. *BioEssays* 38(7):605–612, <https://doi.org/10.1002/bies.201600005>
- [21] Edgar R, Domrachev M, Lash AE (2002) Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res* 30(1):207–210, <https://doi.org/10.1093/nar/30.1.207>
- [22] ENCODE Project Consortium (2011) A user’s guide to the encyclopedia of DNA elements (ENCODE). *PLOS Biol* 9(4):e1001046, <https://doi.org/10.1371/journal.pbio.1001046>
- [23] Frith MC, Hamada M, Horton P (2010) Parameters for accurate genome alignment. *BMC Bioinf* 11:80, <https://doi.org/10.1186/1471-2105-11-80>
- [24] Furey TS (2012) ChIP-seq and beyond: new and improved methodologies to detect and characterize protein-DNA interactions. *Nat Rev Genet* 13(12):840–852, <https://doi.org/10.1038/nrg3306>
- [25] Gupta S, Stamatoyannopoulos JA, Bailey TL, Noble WS (2007) Quantifying similarity between motifs. *Genome Biol* 8(2):R24, <https://doi.org/10.1186/gb-2007-8-2-r24>
- [26] Han J, Back SH, Hur J, Lin YH, Gildersleeve R, Shan J, Yuan CL, Krokowski D, Wang S, Hatzoglou M, Kilberg MS, Sartor MA, Kaufman RJ (2013) ER-stress-induced transcriptional regulation increases protein synthesis leading to cell death. *Nat Cell Biol* 15(5):481–490, <https://doi.org/10.1038/ncb2738>
- [27] Head SR, Komori HK, LaMere SA, Whisenant T, Van Nieuwerburgh F, Salomon DR, Ordoukhanian P (2014) Library construction for next-generation sequencing: overviews and challenges. *BioTechniques* 56(2):61–67, <https://doi.org/10.2144/000114133>



- [28] Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX, Murre C, Singh H, Glass CK (2010) Simple combinations of lineage-determining transcription factors prime *cis*-regulatory elements required for macrophage and B cell identities. *Mol Cell* 38(4):576–589, <https://doi.org/10.1016/j.molcel.2010.05.004>
- [29] Johnson DS, Mortazavi A, Myers RM, Wold B (2007) Genome-wide mapping of *in vivo* protein-DNA interactions. *Science* 316(5830):1497–1502, <https://doi.org/10.1126/science.1141319>
- [30] Joshi S, Davidson G, Le Gras S, Watanabe S, Braun T, Mengus G, Davidson I (2017) TEAD transcription factors are required for normal primary myoblast differentiation *in vitro* and muscle regeneration *in vivo*. *PLOS Genet* 13(2):e1006600, <https://doi.org/10.1371/journal.pgen.1006600>
- [31] Khan A, Mathelier A (2017) Intervene: a tool for intersection and visualization of multiple gene or genomic region sets. *BMC Bioinf* 18:287, <https://doi.org/10.1186/s12859-017-1708-7>
- [32] Kidder BL, Hu G, Zhao K (2011) ChIP-Seq: technical considerations for obtaining high-quality data. *Nat Immunol* 12(10):918–922, <https://doi.org/10.1038/ni.2117>
- [33] King HW, Klose RJ (2017) The pioneer factor OCT4 requires the chromatin remodeller BRG1 to support gene regulatory element function in mouse embryonic stem cells. *eLife* 6:e22631, <https://doi.org/10.7554/eLife.22631>
- [34] Krebs W, Schmidt SV, Goren A, De Nardo D, Labzin L, Bovier A, Ulas T, Theis H, Kraut M, Latz E, Beyer M, Schultze JL (2014) Optimization of transcription factor binding map accuracy utilizing knockout-mouse models. *Nucleic Acids Res* 42(21):13051–13060, <https://doi.org/10.1093/nar/gku1078>
- [35] Krueger F (2012) Trim Galore! [http://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)
- [36] Kulakovskiy I, Levitsky V, Oshchepkov D, Bryzgalov L, Vorontsov I, Makeev V (2013) From binding motifs in ChIP-Seq data to improved models of transcription factor binding sites. *J Bioinf Comput Biol* 11(1):1340004, <https://doi.org/10.1142/S0219720013400040>
- [37] Kulakovskiy IV, Medvedeva YA, Schaefer U, Kasianov AS, Vorontsov IE, Bajic VB, Makeev VJ (2013) HOCOMOCO: a comprehensive collection of human transcription factor binding sites models. *Nucleic Acids Res* 41(Database issue):D195–D202, <https://doi.org/10.1093/nar/gks1089>
- [38] Kundaje A, Jin L, Strattan JS, Maurizio PL (2018) ENCODE transcription factor and histone ChIP-Seq processing pipeline. <https://github.com/ENCODE-DCC/chip-seq-pipeline2>
- [39] Lai E, Clark KL, Burley SK, Darnell JE (1993) Hepatocyte nuclear factor 3/fork head or "winged helix" proteins: a family of transcription factors of diverse biologic function. *Proc Natl Acad Sci USA* 90(22):10421–10423, <https://doi.org/10.1073/pnas.90.22.10421>
- [40] Lambert SA, Jolma A, Campitelli LF, Das PK, Yin Y, Albu M, Chen X, Taipale J, Hughes TR, Weirauch MT (2018) The human transcription factors. *Cell* 172(4):650–665, <https://doi.org/10.1016/j.cell.2018.01.029>
- [41] Lambert SA, Yang AWH, Sasse A, Cowley G, Albu M, Caddick MX, Morris QD, Weirauch MT, Hughes TR (2019) Similarity regression predicts evolution of transcription factor sequence specificity. *Nat Genet* 51(6):981–989, <https://doi.org/10.1038/s41588-019-0411-1>
- [42] Landt SG, Marinov GK, Kundaje A, Kheradpour P, Pauli F, Batzoglu S, Bernstein BE, Bickel P, Brown JB, Cayting P, Chen Y, DeSalvo G, Epstein C, Fisher-Aylor KI, Euskirchen G, Gerstein M, Gertz J, Hartemink AJ, Hoffman MM, Iyer VR, Jung YL, Karmakar S, Kellis M, Kharchenko PV, Li Q, Liu T, Liu XS, Ma L, Milosavljevic A, Myers RM, Park PJ, Pazin MJ, Perry MD, Raha D, Reddy TE, Rozowsky J, Shores N, Sidow A, Slattery M, Stamatoyannopoulos JA, Tolstorukov MY, White KP, Xi S, Farnham PJ, Lieb JD, Wold BJ, Snyder M (2012) ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia. *Genome Res* 22(9):1813–1831, <https://doi.org/10.1101/gr.136184.111>

- [43] Lesluyes T, Johnson J, Machanick P, Bailey TL (2014) Differential motif enrichment analysis of paired ChIP-seq experiments. *BMC Genomics* 15:752, <https://doi.org/10.1186/1471-2164-15-752>
- [44] Lex A, Gehlenborg N, Strobel H, Vuillemot R, Pfister H (2014) UpSet: visualization of intersecting sets. *IEEE T Vis Comput Gr* 20(12):1983–1992, <https://doi.org/10.1109/TVCG.2014.2346248>
- [45] Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25(14):1754–1760, <https://doi.org/10.1093/bioinformatics/btp324>
- [46] Li H, Durbin R (2019) BWA. <https://github.com/lh3/bwa>
- [47] Li Q, Brown JB, Huang H, Bickel PJ (2011) Measuring reproducibility of high-throughput experiments. *Ann Appl Stat* 5(3):1752–1779, <https://doi.org/10.1214/11-AOAS466>
- [48] Lun ATL, Smyth GK (2016) csaw: a Bioconductor package for differential binding analysis of ChIP-seq data using sliding windows. *Nucleic Acids Res* 44(5):e45, <https://doi.org/10.1093/nar/gkv1191>
- [49] Ma W, Noble WS, Bailey TL (2014) Motif-based analysis of large nucleotide data sets using MEME-ChIP. *Nat Protoc* 9(6):1428–1450, <https://doi.org/10.1038/nprot.2014.083>
- [50] Machanick P, Bailey TL (2011) MEME-ChIP: motif analysis of large DNA datasets. *Bioinformatics* 27(12):1696–1697, <https://doi.org/10.1093/bioinformatics/btr189>
- [51] Martin M (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet J* 17(1):10–12, <https://doi.org/10.14806/ej.17.1.200>
- [52] Mathelier A, Fornes O, Arenillas DJ, Chen CY, Denay G, Lee J, Shi W, Shyr C, Tan G, Worsley-Hunt R, Zhang AW, Parcy F, Lenhard B, Sandelin A, Wasserman WW (2016) JASPAR 2016: a major expansion and update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res* 44(D1):D110–D115, <https://doi.org/10.1093/nar/gkv1176>
- [53] McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010) The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* 20(9):1297–1303, <https://doi.org/10.1101/gr.107524.110>
- [54] Merika M, Orkin SH (1993) DNA-binding specificity of GATA family transcription factors. *Mol Cell Biol* 13(7):3999–4010, <https://doi.org/10.1128/mcb.13.7.3999>
- [55] Mitchell PJ, Tjian R (1989) Transcriptional regulation in mammalian cells by sequence-specific DNA binding proteins. *Science* 245(4916):371–378, <https://doi.org/10.1126/science.2667136>
- [56] Newman JRS, Keating AE (2003) Comprehensive identification of human bZIP interactions with coiled-coil arrays. *Science* 300(5628):2097–2101, <https://doi.org/10.1126/science.1084648>
- [57] Park PJ (2009) ChIP-seq: advantages and challenges of a maturing technology. *Nat Rev Genet* 10(10):669–680, <https://doi.org/10.1038/nrg2641>
- [58] Pepke S, Wold B, Mortazavi A (2009) Computation for ChIP-seq and RNA-seq studies. *Nat Methods* 6(11 Suppl):S22–S32, <https://doi.org/10.1038/nmeth.1371>
- [59] Quinlan AR, Hall IM (2010) BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26(6):841–842, <https://doi.org/10.1093/bioinformatics/btq033>
- [60] Robertson G, Hirst M, Bainbridge M, Bilenky M, Zhao Y, Zeng T, Euskirchen G, Bernier B, Varhol R, Delaney A, Thiessen N, Griffith OL, He A, Marra M, Snyder M, Jones S (2007) Genome-wide profiles of STAT1 DNA association using chromatin immunoprecipitation and massively parallel sequencing. *Nat Methods* 4(8):651–657, <https://doi.org/10.1038/nmeth1068>

- [61] Sandelin A, Wasserman WW (2004) Constrained binding site diversity within families of transcription factors enhances pattern discovery bioinformatics. *J Mol Biol* 338(2):207–215, <https://doi.org/10.1016/j.jmb.2004.02.048>
- [62] Savic D, Partridge EC, Newberry KM, Smith SB, Meadows SK, Roberts BS, Mackiewicz M, Mendenhall EM, Myers RM (2015) CETCh-seq: CRISPR epitope tagging ChIP-seq of DNA-binding proteins. *Genome Res* 25(10):1581–1589, <https://doi.org/10.1101/gr.193540.115>
- [63] Schneider TD, Stephens RM (1990) Sequence logos: a new way to display consensus sequences. *Nucleic Acids Res* 18:6097–6100, <https://doi.org/10.1093/nar/18.20.6097>
- [64] Schwenk F, Baron U, Rajewsky K (1995) A *cre*-transgenic mouse strain for the ubiquitous deletion of *loxP*-flanked gene segments including deletion in germ cells. *Nucleic Acids Res* 23(24):5080–5081, <https://doi.org/10.1093/nar/23.24.5080>
- [65] Sternberg N, Hamilton D (1981) Bacteriophage P1 site-specific recombination. *J Mol Biol* 150(4):467–486, [https://doi.org/10.1016/0022-2836\(81\)90375-2](https://doi.org/10.1016/0022-2836(81)90375-2)
- [66] Sullivan AL, Benner C, Heinz S, Huang W, Xie L, Miano JM, Glass CK (2011) Serum response factor utilizes distinct promoter- and enhancer-based mechanisms to regulate cytoskeletal gene expression in macrophages. *Mol Cell Biol* 31(4):861–875, <https://doi.org/10.1128/MCB.00836-10>
- [67] Tarasov A, Vilella AJ, Cuppen E, Nijman IJ, Prins P (2015) Sambamba: fast processing of NGS alignment formats. *Bioinformatics* 31(12):2032–2034, <https://doi.org/10.1093/bioinformatics/btv098>
- [68] Tu S, Shao Z (2017) An introduction to computational tools for differential binding analysis with ChIP-seq data. *Quant Biol* 5(3):226–235, <https://doi.org/10.1007/s40484-017-0111-8>
- [69] Viger RS, Guittot SM, Anttonen M, Wilson DB, Heikinheimo M (2008) Role of the GATA family of transcription factors in endocrine development, function, and disease. *Mol Endocrinol* 22(4):781–798, <https://doi.org/10.1210/me.2007-0513>
- [70] Wei G, Abraham BJ, Yagi R, Jothi R, Cui K, Sharma S, Narlikar L, Northrup DL, Tang Q, Paul WE, Zhu J, Zhao K (2011) Genome-wide analyses of transcription factor GATA3-mediated gene regulation in distinct T cell types. *Immunity* 35(2):299–311, <https://doi.org/10.1016/j.immuni.2011.08.007>
- [71] Weirauch MT, Yang A, Albu M, Cote AG, Montenegro-Montero A, Drewe P, Najafabadi HS, Lambert SA, Mann I, Cook K, Zheng H, Goity A, van Bakel H, Lozano JC, Galli M, Lewsey MG, Huang E, Mukherjee T, Chen X, Reece-Hoyes JS, Govindarajan S, Shaulsky G, Walhout AJM, Bouget FY, Ratsch G, Larrondo LF, Ecker JR, Hughes TR (2014) Determination and inference of eukaryotic transcription factor sequence specificity. *Cell* 158(6):1431–1443, <https://doi.org/10.1016/j.cell.2014.08.009>
- [72] Worsley Hunt R, Wasserman WW (2014) Non-targeted transcription factors motifs are a systemic component of ChIP-seq datasets. *Genome Biol* 15(7):412, <https://doi.org/10.1186/s13059-014-0412-4>
- [73] Zeineddine D, Hammoud AA, Mortada M, Boeuf H (2014) The Oct4 protein: more than a magic stemness marker. *Am J Stem Cells* 3(2):74–82
- [74] Zhang Y, Liu T, Meyer CA, Eeckhoutte J, Johnson DS, Bernstein BE, Nusbaum C, Myers RM, Brown M, Li W, Liu XS (2008) Model-based analysis of ChIP-Seq (MACS). *Genome Biol* 9(9):R137, <https://doi.org/10.1186/gb-2008-9-9-r137>
- [75] Zhang Y, Lin YH, Johnson TD, Rozek LS, Sartor MA (2014) PePr: a peak-calling prioritization pipeline to identify consistent or differential peaks from replicated ChIP-seq data. *Bioinformatics* 30(18):2568–2575, <https://doi.org/10.1093/bioinformatics/btu372>

- [76] Zhao J, Li X, Guo M, Yu J, Yan C (2016) The common stress responsive transcription factor ATF3 binds genomic sites enriched with p300 and H3K27ac for transcriptional regulation. *BMC Genomics* 17(1):335, <https://doi.org/10.1186/s12864-016-2664-8>