

Comparative ChIP-seq (Comp-ChIP-seq): a practical guideline for experimental design and a novel computational methodology

Enrique Blanco¹, Luciano Di Croce^{1,2,3,*} and Sergi Aranda^{1,*}

5

1. Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Dr. Aiguader 88, Barcelona 08003, Spain

2. Universitat Pompeu Fabra (UPF), Barcelona, Spain

3. Institució Catalana de Recerca i Estudis Avançats (ICREA), Pg. Lluís Companys

10 23, Barcelona 08010, Spain

*Correspondence to: Luciano Di Croce (luciano.dicroce@crg.eu) and Sergi Aranda (Sergi.Aranda@crg.eu)

15 **ABSTRACT**

Chromatin immunoprecipitation followed by high-throughput sequencing (ChIP-seq) is a pivotal technique for understanding the functionality of the chromatin-bound factors and for mapping the functional elements of the genome. In order to evaluate cell- and disease-specific changes in the interacting strength of chromatin targets, ChIP-seq signal across multiple conditions must undergo robust normalization. However, this is not possible using the standard ChIP-seq scheme, which lacks a reference for the control of biological and experimental variabilities. While several studies have recently proposed different solutions to circumvent this problem, substantial technical and analytical differences among methodologies could hamper the experimental reproducibility. Here we provide a practical binary decision-making process to experimentally implement a normalizing method for comparative ChIP-seq across different samples. In addition, we evaluate side-by-side the current computational approaches for normalizing using a reference internal genome. Finally, we propose a local regression strategy to accurately normalize ChIP-seq data in a genome-wide manner. Overall, our proposed experimental and computational standard for comparative ChIP-seq (Comp-ChIP-seq) will increase experimental reproducibility, thereby reducing this major confounding factor in interpreting ChIP-seq results.

20
25
30

INTRODUCTION

Chromatin is the macromolecular complex of DNA and histone proteins that packs the genome into its basic structural units of nucleosomes (1). Within chromatin, a plethora of interacting proteins organize the 3D distribution of the genome, regulate multiple gene expression programs, and coordinate the appropriate transmission of genetic and epigenetic information to cellular progeny (1-5). Alterations in the functionality of the proteins associated with chromatin are intimately linked to severe developmental diseases and cancer (6). Indeed, one recent comprehensive analysis of 10,437 cancer exomes has revealed that nearly 40% of all cancer driver genes are well-characterized chromatin associated factors, which include transcription factors, chromatin modifiers, chromatin remodelers, and guardians of genomic stability (7). Due to its biological and pathological relevance, research on chromatin and epigenetics has been a rapidly moving field over the last decade, assisted by the development of novel methods for the high-throughput molecular analysis of the genome.

The development of chromatin immunoprecipitation (ChIP) coupled with the next-generation sequencing (seq) methodologies has been pivotal for characterizing the genomic distribution of a vast collection of chromatin-associated proteins, histone post-translational modifications (PTMs), and histone variants (8-11). The striking impact of the ChIP-seq technology is based in its relative technical simplicity (which allows it to be adopted by most experimental laboratories), its sensitivity and accuracy for mapping the genomic distribution of proteins, and the standardization of the experimental and computational methods to efficiently analyze such a volume of information. Fifteen years ago, the Encyclopedia of DNA Elements (ENCODE) project was launched as a collaborative initiative to catalog the complete set of functional elements in the human genome in selected cell lines (12, 13). More recently, the International Human Epigenome Consortium (IHEC) has generated a comprehensive high-resolution reference map for the epigenome of major primary human cell types (14, 15). The ChIP-seq method has been central for building the cartography of functional elements of the human genome in both of these international collaborative efforts.

Current research efforts aim to identify the changes in the occupancy profiles of chromatin-bound factors or histone modifications in two or more cell types, metabolic states, and/or pathological situations. However, in its traditional scheme, ChIP-seq is essentially a semi-quantitative method that enables the researcher to determine the relative occupancy of one factor in a given genomic region, with respect to the rest of the genome. In other words, there is no direct correlation between the amount of ChIP signal in the output and the biological activity of the binding element in different scenarios (16). Several sources of biological and technical variability can be overlooked, thereby hampering any direct comparisons of ChIP

70 signal strength between different conditions (17). For instance, the method of preparation of the sequencing material by ChIP is a source of technical variability. An apparent increase in genomic occupancy of a chromatin factor could simply be the result of variability in the efficiency of immunoprecipitation or DNA elution between experiments. Moreover, while running the sequencer, a standard practice is to mix equal proportions of barcoded libraries to run the samples in a multiplexed manner. Therefore, even a substantial global reduction of a histone variant occupancy per cell would remain hidden in a ChIP-seq experiment after normalizing by total number of reads (17). Although the consistent replication of ChIP-seq experiments can reveal the biological tendency in the interacting strength of the chromatin factor, a robust normalizing strategy is required to accurately compare ChIP-seq results across experimental conditions.

80 To overcome the influence of technical variabilities in the biological interpretation of ChIP-seq, several groups have reported different strategies based on the use of internal reference controls (spike-in), which provides a feasible solution to accurately normalize comparative ChIP-seq (Table 1). Originally developed to correct gene expression measurement in microarrays and RNA-seq experiments (18, 19), the spike-in strategy is based on combining the experimental sample with an amount of exogenous material (either from another species or synthetically produced) that is constant between experiments. Both the experimental sample and the spike-in are processed and analyzed in parallel. As long as the amount of spike-in ChIP signal is constant, the observable differences in the experimental samples across conditions can be exclusively attributed to biological variation. Eventual differences in the spike-in signal can be computationally equilibrated to eliminate technical variability, and the same correction is then used to normalized the experimental signal.

90 Current strategies for ChIP-seq normalization using spike-in diverge in: *i*) the type of biological material to be used as a reference sample for normalization; *ii*) the capturing method of the spike-in material; and, *iii*) the computational method used to analyze the sequencing data. It is important to mention that the use of different strategies can introduce inconsistencies to the final results, as each method presents its own benefits and limitations. In addition, the existence of several alternatives might complicate decision-making for researchers about when and how to apply a normalization method for comparative ChIP-seq. However, to our knowledge, a benchmarking study involving the existent spike-in alternatives is not available. With the aim of providing a practical, unified reference framework for comparative ChIP-seq analysis, we have now evaluated the benefits and limitations of the experimental methodologies and analytical pipelines currently available. In addition, we have proposed an experimental best practice and, we have designed a bioinformatics pipeline for analyzing comparative ChIP-seq data in a genome-wide manner. If widely used, we believe that this

100

pipeline can serve as a reference for increasing data reproducibility between laboratories, thus overcoming one of the major drawbacks of using ChIP-seq data.

105 **Comparison of current experimental methodologies**

During the last five years, up to four alternative strategies have been proposed based on the spike-in concept to deal with the problem of a lack of comparability between multiple ChIP-seq samples. We have summarized the principal benefits and limitations of each technique in Table 1.

110 In 2014, two different laboratories independently pioneered the development of a similar strategy for comparative ChIP-seq normalization, by introducing xenogenic material from a species different into the experimental model (20, 21). The rationale behind both approaches is that the ChIP-seq signal obtained from a fixed amount of the spike-in material can be used as an internal reference to normalized the experimental ChIP signal across different samples.

115 Thus, Guenther and collaborators mixed *Drosophila melanogaster* S2 cells with human cells before cell lysis (20), while Delorenzi and colleagues added a constant amount of fragmented chromatin from human cells into previously-fragmented mouse chromatin (21). Despite the conceptual similarities, the initial decision on which spike-in material should be added to the experimental sample raises important technical considerations that can influence the resulting

120 biological interpretations. The addition of xenogenic cells of the spike-in material enables the whole procedure to be monitored from the beginning, thereby minimizing the impact of technical variabilities for the biological interpretation. Further, by mixing cells, it is possible to tackle eventual changes in genomic ploidy (e.g. due to genomic instability, or differences in cell cycle progression), thereby providing an estimation of average ChIP-seq signal per cell.

125 This quantitative estimation is not possible when mixing fragmented chromatin. However, the option of mixing cells is only available when the number of cells can be evaluated accurately (e.g. cells growing on a dish), or when the experimental sample and spike-in material are fragmented with the same settings. On the contrary, when number of cells in the sample is uncertain (e.g. animal tissue samples), or when the experimental sample and spike-in material

130 require different settings for fragmentation, the addition of the fragmented spike-in material at chromatin level is a more appropriate option.

In the previous strategies, the spike-in material and the sample are captured using the same antibody, limiting both strategies to using antigens that are highly conserved between the sample and the spike-in material (for instance, histone modifications). To circumvent this

135 problem, Trojer and colleagues introduced a smart solution by using a second antibody for a fly-specific histone variant (H2Av) to capture the spike-in material (22). This strategy aims to avoid the cross-reactivity constraint of the experimental antibody and to reduce any potential variability due to competition between the spike-in control and the experimental material, which usually exceeds the amount of spike-in material by far.

140 In order to overcome both epitope conservation restrictions and the use of a xenogenic spike
material, a fourth normalizing strategy has been recently proposed: rather than using
exogenous material, Guertin and co-workers recommend including a second antibody against
an endogenous target present in the experimental chromatin, as an internal control (23). This
second antibody is used to profile the genomic occupancy of a pervasive chromatin factor (e.g.
145 CTCF) whose genomic distribution is assumed to be unchanged between different cell types
and/or treatments and which is clearly distinguishable from the experimental target (23).
Nonetheless, although this method avoids preparation of xenogenic material, we consider that
it still makes two important assumptions, which could not always be true: i) that the reference
endogenous factor remains stably associated in the different experimental conditions; and, ii)
150 that the genomic distributions of endogenous reference target and the experimental target do
not overlap.

Towards an experimental framework for comparative ChIP-seq experiments

Considering the benefits and limitations of the different strategies (reviewed in Table 1), we
propose the addition of exogenous xenogenic fly cells, whenever possible, and the use of a
155 second antibody against a fly-specific histone variant, as a best practice for comparative ChIP-
seq normalization for mammalian genomes. We recommend using fly material as the spike-in
because: 1) the genome sequence has been extensively assembled; 2) the fly chromatin has
been largely characterized at epigenetic levels; 3) the evolutionary distance between fly and
mammalian genomes is sufficient to allow an unambiguous alignment of the reads (20, 22);
160 and 4) fly cells are relatively easy to culture with standard tissue culture procedures and
instruments. Moreover, the genomic occupancy profile of the fly-specific H2Av is already
characterized (22), which can be extremely useful as an additional control point for assessing
ChIP-seq performance.

Taking into account the previous considerations, we have designed a practical guide under
165 the form of a decision tree to systematically implement a consistent protocol for the addition of
the spike-in material when performing ChIP-seq experiments. In our roadmap, we state some
key questions that are relevant for deciding which type of spike-in to add (*i.e.* fly cells or
fragmented chromatin, diagram of dataflow; Fig. 1). These guidelines take into consideration
that: i) spike-in material should be present in all samples at equal amounts at the earliest step
170 during the ChIP-seq procedure; and (ii) spike-in material should be present in a low-enough
quantity (giving a significantly lower number of reads as that of the experimental reads) to not
interfere with the actual ChIPseq experiment yet still give an accurate normalization in the final
sequencing data. As previously reported, the number of spike-in reads in the final sequencing
step should be at least one million reads, and approximately, 2%–5% of the experimental

175 genome, to minimize the changes in overall material used for ChIP-seq (20, 22). This final
amount of reads can be influenced by the ratio of the mixture as well as by the quality of the
antibody and/or the abundance of the target in the experimental condition. Taking into account
these considerations, and the relative ratio between the size of the fly genome and the two
most widely used mammalian experimental models (mouse and human), we recommend the
180 use of different final mixtures (Fig. 1).

Comparison of current computational methodologies

Traditionally, the standard normalization method uses the total number of mapped reads per
million (RPM) to correct for possible bias introduced by the differences in the sequencing depth
among samples. With the recent spike-in methods (see Table 1), the mapped reads from the
185 spike-in control are used to correct the experimental ChIP-seq signal strength.

At the computational level, the rationale of the different spike-in methodologies is the same:
as long as a constant amount of spike-in material is added to the experimental samples, and
as long as the samples and the spike-in material are processed and analyzed together, the
correction factor computed to eliminate the differences in the spike-in signal can be used to
190 normalize the experimental sample signal. However, the different methodologies differ in the
computational approach to correct the spike-in and, consequently, the experimental sample.

Several authors have proposed to use the number of mapped reads of the spike-in sample
(e.g. *Drosophila melanogaster*) to correct the overall ChIP-seq signal from the main organism
studied (e.g. *human*) (20, 22). Although initially appealing, this fold-change correction presents
195 in our opinion important shortcomings, such as: *i*) the spike-in reads mapped not only along
the ChIP-seq peaks but also over background regions are used for computing the correction
factor; and *ii*) the correction factor is uniformly applied to all experimental reads in the actual
experiment, treating both non-specific and specific signal loci with the same correction value.
Alternatively, the use of a normalization factor computed derived from a pre-defined list of
200 known targets in the spike-in organism was also proposed (21). This option excludes any
confounding background signal, but the correction is only applied to a computational pre-
defined loci list with specific ChIP-signal, thereby relying in the accuracy of the selected
computational peak caller to define the list of positive loci.

Recently, with the aim of overcoming the limitations of the previous approaches, Guertin et al.
205 proposed to apply a linear local regression method (23). This approach computes a correction
coefficient, defined by a linear regression model, for the systematic and gradual correction of
the pre-defined ChIP-seq peaks from the reference. After this, the coefficient is also used to
correct a pre-defined subset of peaks in the experimental sample. In addition, such a linear
correction method implements a statistical approach to calculating the probability and strength

210 of differentially bound loci. Guertin et al. (23) also showed an increased sensitivity (of about 10%) in the detection of differentially bound target loci as compared to the previous absolute fold-change correction. The conceptual improvements of this approach stem from the fact that the correction factor gradually increases along with the informative power (as number of reads) of the peaks. However, the addition of a computational step to pre-select the real signal loci in
215 this strategy, similar to the non-linear correction method, could introduce an additional bias step, as different peak callers compute their own list of peaks. In addition, this analysis would impede the genome-wide evaluation of the signal-to-noise ratio, thereby limiting the informative power of the ChIP-seq.

220 **Benchmarking a novel local regression method for comparative ChIP-seq in a genome-wide manner**

By experimentally implementing the spike-in concept into our routine ChIP-seq experiments, we faced an important limitation while trying to carry out a comparative ChIP-seq in a genome-wide manner. When normalizing an experimental ChIP signal at a genome-wide level, using a constant correction factor from the total number of spike-in reads (absolute fold-change
225 correction), both the positive ChIP signal regions for a target and the background regions were similarly corrected. This introduced uncertainty as to whether the changes in the target occupancy's overall background were biologically meaningful. To overcome this important limitation, we developed a novel method that performs the genome-wide normalization of ChIP-seq data adapting the spike-in control correction to the class of genomic region. Our
230 approach, inspired by the spike-in-based RNA-seq quantification methods described above (18, 19), shares conceptual similarities with the recent linear local correction approach (23) and consists in the application of a local regression, in this case, over all the genome-wide bins determined along the chromosomes (see Methods). Thereby, our method is able to introduce a different correction to each bin in the genome, depending on its class. First, a local
235 regression is computed from the bins in the spike-in genome in order to accommodate the two ChIP-seq conditions compared into the same best-fit line. Next, the values from the real experiment are corrected following the previous local normalization calculated using the spike-in bins (Fig. 2a). Under this approach, the adjustment on a region containing a true ChIP-seq signal will be substantially higher than the change computed for bins located in the
240 background.

To assess the accuracy of our proposal, we compared the performance with the fold-change strategy on a reference dataset (20, 22)(Fig. 2b). We took advantage of the available ChIP-seq data published by Guenther et al. that included fly material as spike-in control (20). In this study, the authors artificially generated a pre-defined ChIP signal gradient for the di-

245 methylation of lysine-79 histone H3 (H3K79me2). To achieve a wide range of distinct conditions, they mixed different proportions of Jurkat cells that had been untreated or treated with a selective inhibitor for the H3K79-methyltransferase DOT1L (EPZ5676). The mixture aims to reflect the global change in the average H3K79me2 level per cell. Finally, they used a constant amount of fly cells as an internal reference control for normalization (Fig. 2b). An appropriate analytical normalization using spike-in should display a quantitative difference between both experimental ChIP signals in the peaks of H3K79me2, while keeping their background levels equilibrated.

We processed the H3K79me2 samples from two completely different conditions: i) the 25:75 (DMSO:EPZ5676) proportion, which has higher levels of H3K79me2; and, ii) the 75:25 proportion, with lower levels of H3K79me2 (Fig. 2b). First, we mapped the resulting sequencing reads to an artificial genome in which we included the human and the fruit fly chromosomes. After separating the mapped reads into human and fly, we segmented the genomes of the sample (human) and the spike-in control fly into bins of 1 Kb. Next, we assigned the maximum absolute ChIP-seq value of H3K79me2 in both conditions for all bins from both genome segmentations (see Methods). These initial values, which were not corrected by any normalization method, were considered to be the raw value (Fig. 2a).

We then used the spike-in data to compute the normalization of each bin using the fold-change (FC) methodology, or our local regression approach (LOESS). For FC correction, we used the total number of aligned fly reads to correct the ChIP signal of both conditions, as previously suggested (20, 22). When running our proposal, we normalized the spike-in sample for the total number of reads, applying the LOESS correction in the fly bins to guide the corrections in the human bins (Fig. 2c). The same correction factors, depending on the density of reads within the bins, were applied to the human bins. The results of both normalization strategies in the experimental ChIP-seq are shown in Fig. 2d. When applying the FC correction in human bins in a genome-wide manner, we quantified a general increase in both in ChIP-signal regions and the background on the 25:75 sample with respect to the 75:25 sample (mean difference of -1.294 vs. -0.548, Fig. 2d). In contrast, when using the LOESS correction, the differences between the 25:75 and the 75:25 samples were still noticeable on bins with positive ChIP signals, while they were balanced on the bins that constitute the background level of the ChIP-seq from the human experiments (-0.789 Vs. -0.086, Fig. 2d). To our knowledge, the distinct gradual normalization from background to positive ChIP signal regions has not previously been considered in any of the existing normalization techniques. We thus believe our proposal presents a technical advance for the genome-wide normalization and for comparison of ChIP-seq experiments.

CONCLUDING REMARKS

The ChIP technique is one of the most widely-used methods in molecular biology (24), since its development over 30 years ago (25-27). The power of this technique has increased dramatically with the advent of the massive parallel sequencing approaches (8-11) and, ChIP-seq experiments have become the universal method to delineate the genome-wide maps of distribution of transcription factors, chromatin remodelers, and histone modifications. Since about a decade ago, the original scheme of ChIP-seq has been maintained substantially unmodified from chromatin isolation and fragmentation, immunoprecipitation using specific antibodies, DNA purification from protein complexes, library preparation, and parallel sequencing. The ChIP-seq experiments endow a large proportion of noise, which can be introduced by the crosslinking artifacts, non-specific binding of the antibody, or the high sensitivity of the parallel sequencing techniques. This significant amount of noise results in many cases in most of the reads mapped into regions of the genome that are unrelated with the chromatin target (i.e. background zones). A precise determination of the signal-to-noise ratio is, therefore, very relevant for determining the occupancy strength of the targets. Limitations due to potential biases introduced by the technical variabilities in the original ChIP-seq scheme have necessitated the development of strategies to implement an internal reference control across samples for further comparisons. However, there are still potential pitfalls in the application of each approach. Thus, this important problem is still open in many aspects. Very recently, an additional spike-in strategy has been developed (28). In this new approach, the authors benefits from the genetic diversity of yeast strains to perform an intra-specie spike-in using different *Saccharomyces cerevisiae* strains as spike-in and experimental specimen. Beyond its utility in lower eukaryotes, this new method exemplifies the aim to meet an experimental need. Considering the different strategies proposed to normalized ChIP-seq data using spike-in, we provide here a major guideline to plan and perform comparative ChIP-seq experiments, and propose an innovative computational approach for the analysis in a genome-wide manner to fill the gap of the experimental need for an accurate interpretation of ChIP-seq data.

Our proposed standard for the genome-wide comparative ChIP-seq signal has shown to be very effective for the precise comparison of ChIP signals across samples without a pre-defined selection of the loci. This approach is able to correct for possible technical bias and to compute a local correction factor, thereby minimizing the impact of the correction over non-occupied genomic regions. We strongly believe that the systematic use of spike-in references in the ChIP-seq experiments will provide a more precise picture of the dynamics of epigenome in different conditions.

The epigenetic research field is systematically searching to increase the sensitivity for minimizing the amount of cells required for analysis, with the aim of reaching a confident single-cell ChIP-seq. The high variability produced by manipulating single-cell events shall benefit from the introduction of internal reference controls, which are at the conceptual bottom line of the methods reported here. We expect that single-cell Comp-ChIP-seq will bring this technique into a third revolution for tracking the genomic effects of molecules bound to single loci.

METHODS

First, a genome index was generated that combined the human chromosomes (assembly: hg19) and the fruit fly chromosomes (assembly: dm3). The fly chromosomes were labeled using the “_FLY” tag. Additionally, the human and fruit fly gene transcripts, as annotated by RefSeq (29), were merged into a single catalog of genes. The raw data corresponding to the samples Jurkat_K79_25%_R1 (GSM1465005) and Jurkat_K79_75%_R1 (GSM1465007) from (20) was downloaded from the GEO record GSE60104. BOWTIE (30) was run to map both samples of reads over the human+fly genome index (previously generated). By using the “_FLY” tag, the mapped reads corresponding to the spike-in control were then separated from the human experimental values. To perform the genome-wide normalization benchmarking, we segmented the human and fly chromosomes into bins of 1 Kb. Next, we assigned the highest value of every ChIP-seq profile of read counts inside each bin. The following transformations were necessary to perform the regression plots: *i*) absolute values, in millions of reads, which were used as raw values; *ii*) absolute values divided by the total number of human+fly reads per sample, for the traditional normalization; and *iii*) absolute values divided by the total number of fly reads per sample, for the fold-change normalization. The LOESS function from the R library `affy` was applied to the traditional normalization values, to perform the local regression of data, inspired in a similar treatment proposed for RNA-seq normalization of the RPKMs of spike-in controls (18). We instructed the `loess` function to use the adjustment on the values in the spike-in genome as a subset to guide the normalization of the human values. MACS (31) was used to identify the list of ChIPseq peaks along both genomes. To discriminate between bins that contain ChIPseq peaks and bins that constitute the background, the overlap was calculated between MACS peaks and the coordinates of the segmentation bins at human and fly chromosomes. We used boxplots to represent the distribution of normalized values that belongs to each class of bins.

350 **REFERENCES**

1. P. Cramer, A tale of chromatin and transcription in 100 structures. *Cell* **159**, 985-994 (2014).
2. B. Bonev, G. Cavalli, Organization and function of the 3D genome. *Nat Rev Genet* **17**, 661-678 (2016).
- 355 3. G. Almouzni, H. Cedar, Maintenance of Epigenetic Information. *Cold Spring Harb Perspect Biol* **8**, (2016).
4. D. M. Gilbert *et al.*, Space and time in the nucleus: developmental control of replication timing and chromosome architecture. *Cold Spring Harb Symp Quant Biol* **75**, 143-153 (2010).
- 360 5. S. Aranda, G. Mas, L. Di Croce, Regulation of gene transcription by Polycomb proteins. *Sci Adv* **1**, e1500737 (2015).
6. A. C. Mirabella, B. M. Foster, T. Bartke, Chromatin deregulation in disease. *Chromosoma* **125**, 75-93 (2016).
- 365 7. M. H. Bailey *et al.*, Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell* **173**, 371-385 e318 (2018).
8. I. Albert *et al.*, Translational and rotational settings of H2A.Z nucleosomes across the *Saccharomyces cerevisiae* genome. *Nature* **446**, 572-576 (2007).
9. A. Barski *et al.*, High-resolution profiling of histone methylations in the human genome. *Cell* **129**, 823-837 (2007).
- 370 10. D. S. Johnson, A. Mortazavi, R. M. Myers, B. Wold, Genome-wide mapping of in vivo protein-DNA interactions. *Science* **316**, 1497-1502 (2007).
11. T. S. Mikkelsen *et al.*, Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature* **448**, 553-560 (2007).
- 375 12. E. P. Consortium, The ENCODE (ENCyclopedia Of DNA Elements) Project. *Science* **306**, 636-640 (2004).
13. E. P. Consortium, An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74 (2012).
- 380 14. H. G. Stunnenberg, C. International Human Epigenome, M. Hirst, The International Human Epigenome Consortium: A Blueprint for Scientific Collaboration and Discovery. *Cell* **167**, 1145-1149 (2016).
15. M. Skipper *et al.*, Presenting the epigenome roadmap. *Nature* **518**, 313 (2015).
16. S. G. Landt *et al.*, ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia. *Genome Res* **22**, 1813-1831 (2012).
- 385 17. K. Chen *et al.*, The Overlooked Fact: Fundamental Need for Spike-In Control for Virtually All Genome-Wide Analyses. *Mol Cell Biol* **36**, 662-667 (2015).
18. J. Loven *et al.*, Revisiting global gene expression analysis. *Cell* **151**, 476-482 (2012).

19. F. Taruttis *et al.*, External calibration with *Drosophila* whole-cell spike-ins delivers absolute mRNA fold changes from human RNA-Seq and qPCR data. *Biotechniques* **62**, 53-61 (2017).
- 390 20. D. A. Orlando *et al.*, Quantitative ChIP-Seq normalization reveals global modulation of the epigenome. *Cell reports* **9**, 1163-1170 (2014).
21. N. Bonhoure *et al.*, Quantifying ChIP-seq data: a spiking method providing an internal reference for sample-to-sample normalization. *Genome Res* **24**, 1157-1168 (2014).
- 395 22. B. Egan *et al.*, An Alternative Approach to ChIP-Seq Normalization Enables Detection of Genome-Wide Changes in Histone H3 Lysine 27 Trimethylation upon EZH2 Inhibition. *PloS one* **11**, e0166438 (2016).
23. M. J. Guertin, A. E. Cullen, F. Markowitz, A. N. Holding, Parallel factor ChIP provides essential internal control for quantitative differential ChIP-seq. *Nucleic Acids Res* **46**, e75 (2018).
- 400 24. S. Aranda, Y. Shi, L. Di Croce, Chromatin and Epigenetics at the Forefront: Finding Clues among Peaks. *Mol Cell Biol* **36**, 2432-2439 (2016).
25. M. J. Solomon, P. L. Larsen, A. Varshavsky, Mapping protein-DNA interactions in vivo with formaldehyde: evidence that histone H4 is retained on a highly transcribed gene. *Cell* **53**, 937-947 (1988).
- 405 26. D. S. Gilmour, J. T. Lis, Detecting protein-DNA interactions in vivo: distribution of RNA polymerase on specific bacterial genes. *Proc Natl Acad Sci U S A* **81**, 4275-4279 (1984).
27. D. S. Gilmour, J. T. Lis, In vivo interactions of RNA polymerase II with genes of *Drosophila melanogaster*. *Mol Cell Biol* **5**, 2009-2018 (1985).
- 410 28. L. A. Vale-Silva, T. E. Markowitz, A. Hochwagen, SNP-ChIP: a versatile and tag-free method to quantify changes in protein binding across the genome. *BMC Genomics* **20**, 54 (2019).
29. N. A. O'Leary *et al.*, Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**, D733-745 (2016).
- 415 30. B. Langmead, C. Trapnell, M. Pop, S. L. Salzberg, Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* **10**, R25 (2009).
- 420 31. Y. Zhang *et al.*, Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137 (2008).

ACKNOWLEDGMENTS

We specially thank Dr. Cecilia Ballare and Dr. Pedro Vizan, as well as all the
425 members of the Di Croce laboratory for critical reading of the manuscript and
insightful discussions and V.A. Raker for scientific editing. **Funding:** We
acknowledge support from the Spanish Ministry of Economy, Industry and
Competitiveness to the EMBL partnership, Centro de Excelencia Severo Ochoa, the
430 CERCA Programme / Generalitat de Catalunya. The Secretary for Universities and
Research of the Ministry of Economy and Knowledge of the Government of
Catalonia and the Lady Tata Memorial Trust (to S.A.). The Spanish Ministerio de
Educación y Ciencia (BFU2016-75008-P), AGAUR, and La Marato TV3 (to L.D.C).
Author contributions: E.B. conceived and performed the bioinformatics analysis of
435 deep sequencing data and contributed to writing the manuscript. L.D.C. contributed
to data analysis and interpretation and to writing the manuscript. S.A. conceived and
planned this project, perform data analysis and interpretation, and wrote the
manuscript with input from the coauthors. **Competing interests:** The authors
declare that they have no competing interests.

440 **FIGURE LEGENDS**

TABLE 1: Table summarizing the different methodologies developed to normalized ChIP-seq signal with internal controls.

For each spike-in method we initially provide the following basic information: bibliographical
445 reference, type of material used to normalize and the method used to capture the reference
control. Next, we assess the strengths and weaknesses of each method according to several
parameters: cell counting requirements, the constrain in epitope conservation, the capacity to
estimate the average ChIP signal per cell, the possibility to monitor the fragmentation efficiency
in the reference material and the stability of the ChIP-seq signal in the reference control.
450 Finally, we indicate the computational method undertaken for each analysis.

Figure 1: Flow chart to select the most appropriated strategy for a Comp-ChIP-seq experiment.

Figure 2: Novel computational approach for normalizing ChIP-seq data using spike-in ChIP-signal in a genome-wide manner. **a** A diagram summarizing the computational analysis protocol (see Methods section for details). **b** Scheme representing the experimental approach undertaken in (20) to generate a gradual ChIP-seq signal for H3K79me2. **c** Scatterplots showing the distribution of fly bins in both conditions accordingly to the maximum value of H3K79me2 on each bin. From top to bottom, we show the raw values, the values
460 adjusted by the sequencing depth and the final values corrected by the resulting local regression by LOESS to the best-fit line. **d** Box plots representing the distribution of H3K79me2 ChIP-seq signal after absolute fold-change normalization (FC), or LOESS normalization, using the spike-in material. The number below each pair of distributions corresponds to the differential of the mean between conditions resulting from each normalization. We show the
465 values of each experiment distinguishing between bins overlapping H3K79me2 peaks and bins over the background regions.

Table 1

Name	ChIP-Rx	Single-antibody ChIP	Two-antibodies ChIP	Parallel-factor ChIP
Reference	<i>Orlando, D.A. et al, 2014</i>	<i>Bonhoure, N et al. 2014</i>	<i>Egan, B. et al, 2016</i>	<i>Guertin, M.J. et al, 2018</i>
Normalizing sample	Xenogenic cells	Xenogenic chromatin	Xenogenic chromatin	Parallel ChIP with sample material
Capturing method of spike-in material	Antigen conservation between sample and spike-in material	Antigen conservation between sample and spike-in material	Specific antibody for spike-in material (H2Av)	Second antibody for endogenous protein
Cell counting	Required	Not required	Not required	Not required
Epitope conservation	Required	Required	Not required	Not required
Estimation average ChIP signal per cells	Possible	Not possible	Not possible	Not possible
Monitoring fragmentation on spike-in material	Not possible	Possible	Possible	Possible
ChIP -seq profiling of spikein material	Constant	Constant	Constant	Assumed to be constant
Computational analysis	Absolute fold-change correction	Signal fold-change correction	Absolute fold-change correction	Linear local correction

Figure 1

Guideline for spike-in addition

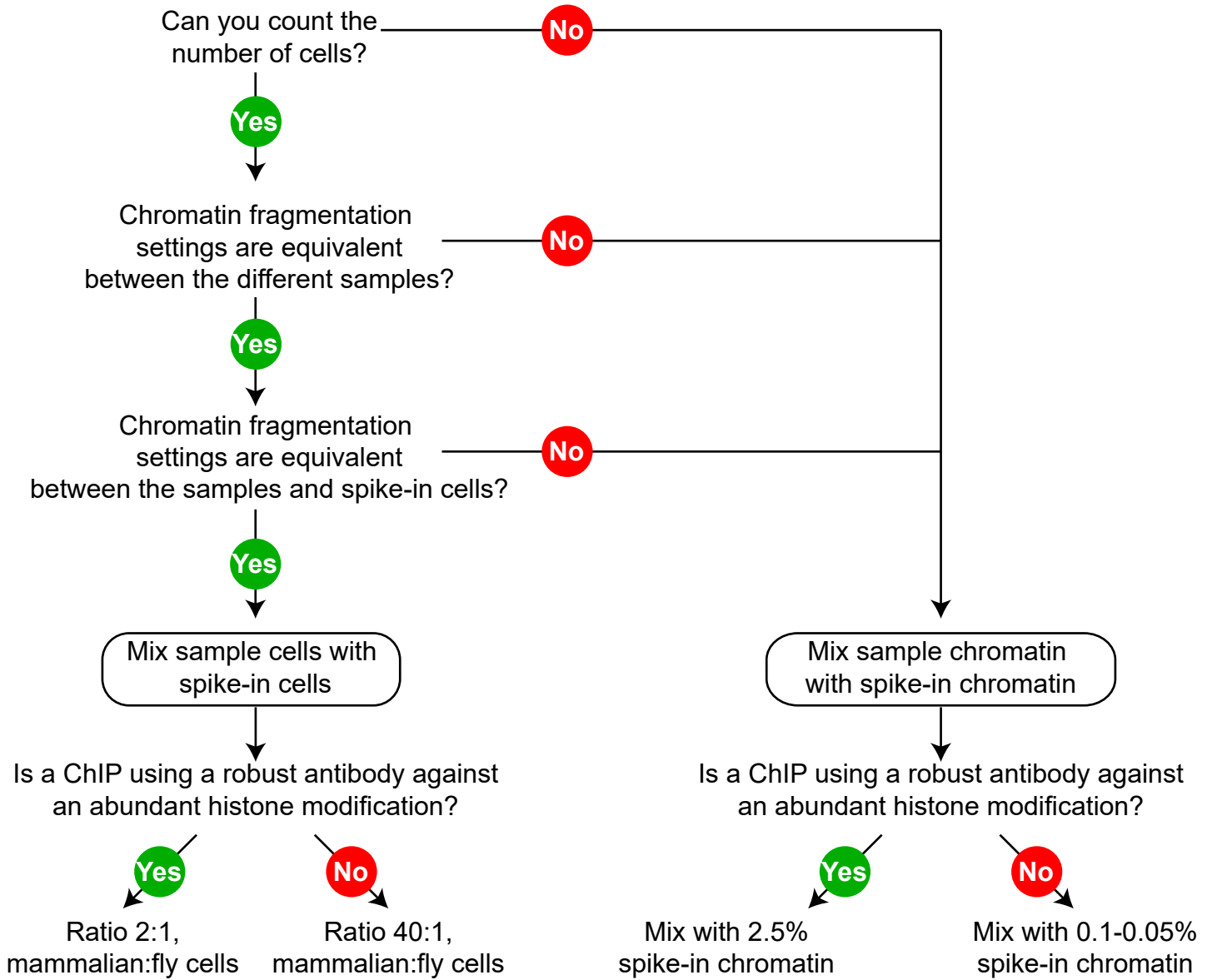
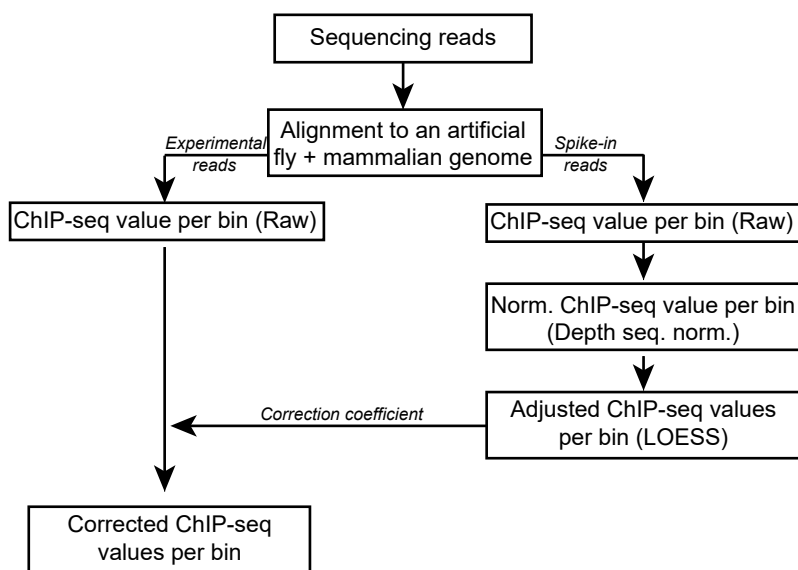
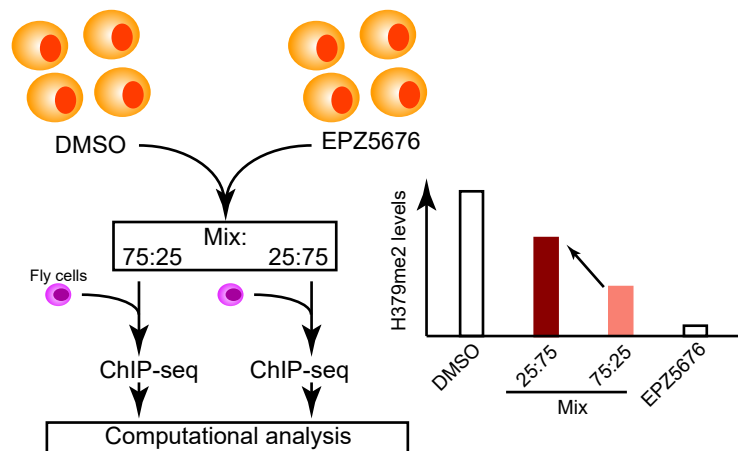


Figure 2

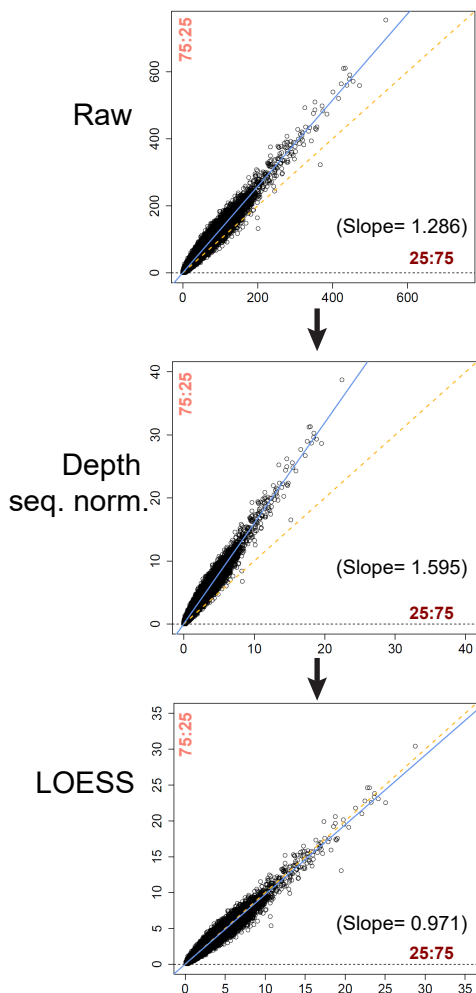
A Computational pipeline for genome-wide LOESS correction with spike-in data



B Data-set; Orlando et al. 2014



C ChIP-seq values on 120,397 fly genomic bins (1Kb)



D

