



**Citation:** Rana V, Peng J, Pan C, Lyu H, Cheng A, Kim M, et al. (2024) Interpretable online network dictionary learning for inferring long-range chromatin interactions. PLoS Comput Biol 20(5): e1012095. https://doi.org/10.1371/journal.pcbi.1012095

**Editor:** Tamar Schlick, New York University, UNITED STATES

**Received:** December 16, 2023 **Accepted:** April 20, 2024

Published: May 16, 2024

Copyright: © 2024 Rana et al. This is an open access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The code is available at https://github.com/rana95vishal/chromatin\_DL The complete dataset is available at https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE109355.

**Funding:** The work was supported by the National Science Foundation grants #1956384 (AC, MK, and OM), #2206296 (OM) and grant CZI DAF 2022-249217 (OM and MK). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

RESEARCH ARTICLE

# Interpretable online network dictionary learning for inferring long-range chromatin interactions

Vishal Rana<sup>1</sup>, Jianhao Peng<sup>1</sup>, Chao Pan<sup>1</sup>, Hanbaek Lyu<sup>2</sup>, Albert Cheng<sup>3</sup>, Minji Kim<sup>4</sup>, Olgica Milenkovic<sub>6</sub><sup>1</sup>\*

1 Department of Electrical and Computer Engineering, University of Illinois, Urbana-Champaign, Illinois, United States of America, 2 Department of Mathematics, University of Wisconsin - Madison, Madison, Wisconsin, United States of America, 3 School of Biological and Health Systems Engineering, Arizona State University, Phoenix, Arizona, United States of America, 4 Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, Michigan, United States of America

\* milenkov@illinois.edu

# **Abstract**

Dictionary learning (DL), implemented via matrix factorization (MF), is commonly used in computational biology to tackle ubiquitous clustering problems. The method is favored due to its conceptual simplicity and relatively low computational complexity. However, DL algorithms produce results that lack interpretability in terms of real biological data. Additionally, they are not optimized for graph-structured data and hence often fail to handle them in a scalable manner.

In order to address these limitations, we propose a novel DL algorithm called *online convex network dictionary learning* (online cvxNDL). Unlike classical DL algorithms, online cvxNDL is implemented via MF and designed to handle extremely large datasets by virtue of its online nature. Importantly, it enables the interpretation of dictionary elements, which serve as cluster representatives, through convex combinations of real measurements. Moreover, the algorithm can be applied to data with a network structure by incorporating specialized subnetwork sampling techniques.

To demonstrate the utility of our approach, we apply cvxNDL on 3D-genome RNAPII ChIA-Drop data with the goal of identifying important long-range interaction patterns (long-range dictionary elements). ChIA-Drop probes higher-order interactions, and produces data in the form of hypergraphs whose nodes represent genomic fragments. The hyperedges represent observed physical contacts. Our hypergraph model analysis has the objective of creating an interpretable dictionary of long-range interaction patterns that accurately represent global chromatin physical contact maps. Through the use of dictionary information, one can also associate the contact maps with RNA transcripts and infer cellular functions.

To accomplish the task at hand, we focus on RNAPII-enriched ChIA-Drop data from *Drosophila Melanogaster* S2 cell lines. Our results offer two key insights. First, we demonstrate that online cvxNDL retains the accuracy of classical DL (MF) methods while simultaneously ensuring unique interpretability and scalability. Second, we identify distinct collections of proximal and distal interaction patterns involving chromatin elements shared by related

**Competing interests:** The authors have declared that no competing interests exist.

processes across different chromosomes, as well as patterns unique to specific chromosomes. To associate the dictionary elements with biological properties of the corresponding chromatin regions, we employ Gene Ontology (GO) enrichment analysis and perform multiple RNA coexpression studies.

# Author summary

We introduce a novel method for dictionary learning termed *online convex Network Dictionary Learning* (online cvxNDL). The method operates in an online manner and utilizes representative subnetworks of a network dataset as dictionary elements. A key feature of online cvxNDL is its ability to work with graph-structured data and generate dictionary elements that represent convex combinations of real data points, thus ensuring interpretability.

Online cvxNDL is used to investigate long-range chromatin interactions in S2 cell lines of *Drosophila Melanogaster* obtained through RNAPII ChIA-Drop measurements represented as hypergraphs. The results show that dictionary elements can accurately and efficiently reconstruct the original interactions present in the data, even when subjected to convexity constraints. To shed light on the biological relevance of the identified dictionaries, we perform Gene Ontology enrichment and RNA-seq coexpression analyses. These studies uncover multiple long-range interaction patterns that are chromosome-specific. Furthermore, the findings affirm the significance of convex dictionaries in representing TADs cross-validated by imaging methods (such as 3-color FISH (fluorescence in situ hybridization)).

### Introduction

Dictionary learning (DL) is a widely used method in learning and computational biology for approximating a matrix through sparse linear combinations of dictionary elements. DL has been used in various applications such as clustering, denoising, data compression, and extracting low-dimensional patterns [1–8]. For example, DL is used to cluster data points since dictionary elements essentially represent centroids of clusters. DL can perform denoising by combining only the highest-score dictionary elements to reconstruct the input; in this case, the low-score dictionary elements reflect the distortion in the data due to noise. DL can also perform efficient data compression by storing only the dictionary elements and associated weights needed for reconstruction. In addition, DL can be used to extract low-dimensional patterns from complex high-dimensional inputs.

However, standard DL methods [9, 10] suffer from interpretability and scalability issues and are primarily applied to *unstructured* data. To address interpretability issues for unstructured data, convex matrix factorization was introduced in [11]. Convex matrix factorization requires that the dictionary elements be convex combinations of real data points, thereby introducing a constraint that adds to the computational complexity of the method. At the same time, to improve scalability, DL and convex DL algorithms can be adapted to online settings [12, 13]. Network DL (NDL), introduced in [14], operates on graph-structured data and samples subnetworks via Markov Chain Monte Carlo (MCMC) methods [14–16] to efficiently and accurately identify a small number of subnetwork dictionary elements that best explain subgraph-level interactions of the entire global network. These dictionary elements

learned by the original NDL algorithm only provide 'latent' subgraph structures that are not necessarily associated with specific subgraphs in the network. When applied to gene interaction networks, such latent subnetworks cannot be associated with specific genomic regions or viewed as physical interactions between genomic loci, making the method biologically uninterpretable.

To address the shortcoming of online NDL, we propose online cvxNDL, a novel NDL method that combines the MCMC sampling technique from [14] with convexity constraints on the matrix representation of sampled subnetworks. These constraints are handled through the concept of "dictionary element representatives," which are essentially adjacency matrices of real subnetworks of the input network. The representatives are used as building blocks of actual dictionary elements. More precisely, dictionary elements are convex combinations of small subsets of representatives. This allows us to map the dictionary element entries to actual genomic regions and view them as real physical interactions. The online learning component is handled via sequential updates of the best choice of representative elements, complementing the approach proposed in [13] for unstructured data. This formulation ensures interpretability of the results and allows for scaling to large datasets.

The utility of online cvxNDL is demonstrated by performing an extensive analysis of 3D chromatin interaction data generated by the RNAPII ChIA-Drop [17] technique. Chromatin 3D structures play a crucial role in gene regulation [18, 19] and have traditionally been measured using "bulk" sequencing methods, such as Hi-C [20] and ChIA-PET [21, 22]. However, due to the proximity ligation step, these methods can only capture pairwise contacts and fail to extract potential multiway interactions that exist in the cell. Further, these methods operate on a population of millions of molecules and therefore only provide information about population averages. ChIA-Drop, by contrast, mitigates these issues by employing droplet-based barcode-linked sequencing to capture multiway chromatin interactions at the single-molecule level, enabling the detection of short- and long-range interactions involving multiple genomic loci. Note that, more specifically, RNAPII ChIA-Drop data elucidates interactions among regulatory elements such as enhancers and promoters, which warrants contrasting/combining it with RNA-seq data.

The cvxNDL method is first tested on synthetic data, and, subsequently, on real-world RNAPII ChIA-Drop data pertaining to chromosomes of *Drosophila Melanogaster* Schneider 2 (S2) phagocytic cell lines (Due to the limited number of complete ChIA-Drop datasets, we only report findings for cell-lines also studied in [17]). For simplicity, we henceforth refer to the latter as ChIA-Drop data (Our method is designed to handle multiway interactions generated by ChIA-Drop experiments and to generate dictionary elements that capture fundamental chromatin interactions. However, it can also be directly applied to other conformation maps, including Hi-C matrices, but without the hypergraph preprocessing steps). Our findings are multi-fold.

First, we provide dictionary elements that can be used to represent chromatin interactions in a succinct and highly accurate manner.

Second, we discover significant differences between the long-range interactions captured by dictionary elements of different chromosomes. These differences can also be summarized via the average distance between interacting genomic loci and the densities of interactions.

Third, we perform Gene Ontology (GO) enrichment analysis to gain insights into the collective functionality of the genomic regions represented by the dictionary elements of different chromosomes. As an example, for chromosomes 2L and 2R, our GO enrichment analysis reveals significant enrichment in several important terms related to reproduction, oocyte differentiation, and embryonic development. Likewise, chromosomes 3L and 3R are enriched in key GO terms associated with blood circulation and response to heat and cold.

Fourth, to further validate the utility of the dictionary elements, we perform an RNA-Seq coexpression analysis using data from independent experiments conducted on *Drosophila Melanogaster* S2 cell lines, available through the NCBI Sequence Read Archive [23]. We show that genes associated with a given dictionary element exhibit high levels of coexpression, as validated on TAD interactions T1-T4 and R1-R4 [17]. Notably, a small subset of our dictionary elements is able to accurately represent these TAD regions and their multiway interactions, confirming the capability of our method to effectively capture complex patterns of both short-and long-range interactions. In addition, we map our dictionary elements onto interaction networks, including the STRING protein-protein interaction network [24], as well as large gene expression repositories like FlyMine. We observe closely coordinated coexpression among the identified genes, further supporting the biological relevance of the identified dictionary elements.

With its unique features, our new interpretable method for dictionary learning adds to the growing literature on machine learning approaches that aim to elucidate properties of chromatin interactions [25–28].

### Results and discussion

We first provide an intuitive, high-level overview of the steps of the interpretable dictionary learning method, as illustrated in Fig 1. The figure describes the most important global ideas behind our novel online cvxNDL pipeline. A rigorous mathematical formulation of the problem and relevant analyses are delegated to the Methods Section, while detailed algorithmic methods are available in Section B in S1 Text.

Chromatin interactions are commonly represented as contact maps. A contact map can be viewed as a hypergraph, where nodes represent genomic loci and two or more such nodes are connected through hyperedges to represent experimentally observed multiway chromatin interactions. Since it is challenging to work with hypergraphs directly, the first step is to transform a hypergraph into an ordinary network (graph), which we tacitly assume is connected. For this purpose, we employ *clique expansion* [29, 30], as shown in Fig 1B. Clique expansion converts a hyperedge into a clique (a fully connected network) and therefore preserves all interactions encapsulated by the hyperedge. However, large hyperedges covering roughly 10 or more nodes in the network can introduce distortion by creating new cliques that do not correspond to any multiway interaction, as shown in Fig 1C [31]. The frequency of such large hyperedges and the total number of hyperedges in chromatin interaction data is limited (i.e., the hypergraph is sparse, see Table A in S1 Text). This renders the distortion due to the hypergraph-to-network conversion process negligible.

To generate an online sample from the clique-expanded input network, we use a subnetwork sampling procedure shown in Fig 1D. We consider a small template network consisting of a fixed number of nodes and search for induced subnetworks in the input that contain the template network topology. These induced subnetworks can be rigorously characterized via *homomorphisms* and are discussed in detail in the Methods Section. An example of a homomorphism is shown in Fig 1D. Throughout our analysis, we will *exclusively focus on path homomorphisms* because they are most suitable for the biological problem investigated. To generate a sequence of online samples from the input network, we employ MCMC sampling. Given a path sample at discrete time t, the next sample at time t+1 is generated by selecting a new node uniformly at random from the neighborhood of the sample at time t and calculating its probability of acceptance  $\beta$ , explained in the Methods Section. If this new node is accepted, we perform a *directed* random walk starting at the selected node, otherwise, we restart the random walk from the first node of the sample at time t. Note that the input network is undirected

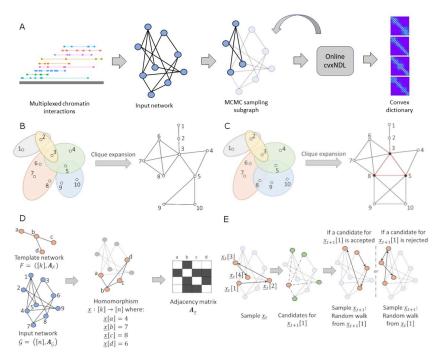
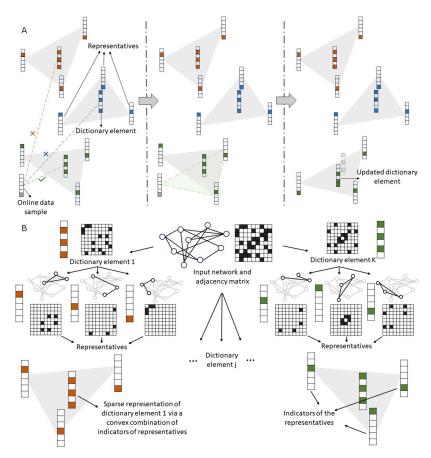


Fig 1. (a) Workflow of the dictionary learning method. Multiway (multiplexed) chromatin interactions represented as hyperedges are clique expanded into standard networks and combined to create input networks for the algorithm. MCMC subnetwork sampling is then used to generate samples for initialization and online updates during iterative optimization of the objective function, resulting in convex dictionary elements. (b) Illustration of the clique expansion process. Hyperedges are subsets of indexed nodes shaded with the same color. (c) Illustration of clique expansion distortion. There is no hyperedge including nodes 3, 5, and 8 (colored red), and this 3-clique only exists due to shared nodes/edges of "real" hyperedges. Such distortion is negligible when the number of large hyperedges is limited. (d) Subnetwork sampling and the notion of a motif homomorphism. These correspond to subnetworks of the input network induced by a fixed number of nodes that contain a template motif topology. The set of homomorphisms  $\operatorname{Hom}(F,\mathcal{G})$  for a network  $\mathcal{G}$  and the template network F are defined in the Methods Section (Eq.7). Also depicted are an example homomorphism  $\underline{x} \in \text{Hom}(F, \mathcal{G})$  and its induced adjacency matrix  $\mathbf{A}_{x}$  for an input network  $\mathcal{G}$  with 9 nodes. The template F is a star network on 4 nodes. In the adjacency matrix, a black field indicates 1, while a white field indicates 0. (e) Workflow of the MCMC sampling algorithm for path homomorphisms. Given a sample  $\underline{x}$ , at time t, obtained via a directed random walk from an initial state in the input network,  $\underline{x}_t[1]$ , we generate a sample  $\underline{x}_t+1$  at time t+1 by choosing uniformly at random a node  $\nu$  from the neighborhood of  $x_{i}[1]$  (marked in green) and calculating a probability of acceptance  $\beta$ . If node  $\nu$  is accepted, we initiate a new directed random walk from  $\nu$ , otherwise, we restart a directed random walk from  $\underline{x}_t[1]$  to generate a new sample.

while only the sampling method requires a directed walk as the order of the labeled nodes matters. (see Fig 1E).

MCMC sampling is used to generate a sequence of samples to initialize a dictionary with *K* dictionary elements, where *K* is chosen based on the properties of the dataset. Each of the dictionary elements is represented as a convex combination of a *small* (sparse) set of *representatives* that are real biological observations. The convex hull of these representatives is termed the *representative region* of the dictionary element. As a result, the vertices of the representative regions comprise a collection of MCMC-generated real-world samples. Fig 2A shows the organization of a dictionary as a collection of dictionary elements, representatives, and representative regions.

After initialization, we perform iterative optimization of the DL objective function using online samples, again generated via the MCMC method. More precisely, at each iteration, we compute the distance between the new sample and every current estimate of dictionary



**Fig 2.** (a) Organization of a dictionary comprising K dictionary elements that are convex combinations of real representative subnetworks. Each dictionary element itself is a sparse *convex combination* of a set of representatives which are small subnetworks of the input real-world network. In the example, there are 6 options for the representatives, and inclusion of a representative into a dictionary element is indicated by a colored entry in a 6-dimensional indicator column-vector. Each of the 6 representatives corresponds to a subnetwork of the input network with a fixed number of nodes (3 for our example). The dictionary element is generated by a convex combination of the corresponding adjacency matrices of its corresponding representative subnetworks. For the example, the resulting dictionary elements are  $9 \times 9$  matrices. (b) Illustration of the representative region update. When an online data sample is observed, the distance of the sample to each of the current dictionary elements is computed and the sample is assigned to the representative region of the nearest dictionary element. From this expanded set of representatives, one representative is carefully selected for removal to improve the objective. The new dictionary element is then obtained as an optimized convex combination of the updated set of representatives.

elements. Subsequently, we assign the sample to the representative region of the nearest dictionary element, which leads to an increase in the size of the set of representatives associated with the dictionary element. From this expanded set of representatives, we carefully select one representative for removal, maximizing the improvement in the quality of our dictionary element and the objective function. It is possible that the removed representative is the newly added data sample assigned to the representative region. In this case, the dictionary element remains unchanged. Otherwise, it is obtained as a convex combination of the updated set of representatives. After observing sufficiently many online samples, the algorithm converges to an accurate set of dictionary elements or the procedure terminates without convergence (in which case we declare a failure and restart the learning process). In our experiments, we never terminated with failure, but due to the lack of provable convergence guarantees for real-world datasets, such scenarios cannot be precluded. The update procedure is shown in Fig 2B.

**Fig 3. Generation of ChIA-Drop data.** ChIA-Drop [17] adopts a droplet-based barcode-linked technique to reveal multiway chromatin interactions at a single molecule level. Chromatin samples are crosslinked and fragmented without a proximity ligation step. The samples are enriched for informative fragments through antibody pull-down.

We applied the method outlined above to RNAPII-enriched ChIA-Drop data from *Drosophila Melanogaster* S2 cells, using a dm3 reference genome [17], to learn dictionaries of chromatin interactions. Fig 3 provides an illustration of the ChIA-Drop pipeline.

We preprocessed the RNAPII ChIA-Drop data to remove fragments mapped to the repetitive regions in the genome and performed an MIA-Sig enrichment test with FDR 0.1 [32]. Only the hyperedges that passed this test were used in subsequent analysis. The highest interaction resolution of our method is dictated by the technology used to generate the data. Since ChIA-Drop experiments involved genomic fragments of length  $\leq$  630 bases [17], we binned chromosomal genetic sequences into fragments of 500 bases each and used the midpoint of each fragment for distance evaluations and dictionary element mappings onto chromatin order. These bins of 500 consecutive bases form the nodes of the hypergraph for each chromosome, while the set of filtered multiway interactions form the hyperedges. The dataset hence includes 45, 938, 42, 292, 49, 072, and 55, 795 nodes and 36, 140, 28, 387, 53, 006, 45, 530 hyperedges for chromosome chr2L, chr2R, chr3L and chr3R respectively. The distribution of the hyperedge sizes is given in Table A in S1 Text. To create networks from hypergraphs, we converted the multiway interactions into cliques. The clique-expanded input network has 113, 606, 85, 316, 161, 590, and 143, 370 edges respectively. Although the ChIA-Drop data comprises interactions from six chromosomes chr2L, chr2R, chr3L, chr3R, chr4 and chrX, since chr4 and chrX are relatively short regions and most of the functional genes are located on chr2L, chr2R, chr3L, and chr3R, we focus our experiments only on the latter.

In the analyses, we fix the number of dictionary elements to K = 25. Clearly, other genomic datasets may benefit from a different choice of the parameter K, which has to be fine-tuned for each different dataset. Also, as template subnetworks, we use paths, since paths are the simplest and most common network motifs, especially in chromatin interaction data (most contact measurements are proximal due to the linear chromosome order). We select paths of length 21 nodes (i.e.,  $21 \times 500$  bases). Once again, both the choice of the subnetwork (motif) and its number of constituent nodes is data dependent. The detailed explanation below justifies our parameter choices for the *Drosophila* dataset.

The typical range of long-range interactions in chromatin structures depends on the species/reference genome. For *Drosophila Melanogaster*, TADs are 10, 000–100, 000 bases long, while loops are usually (much) shorter than 10, 000 bases [17, 33]. This suggests using 10, 000 bases as an approximate lower bound for the length of long-range interactions. In addition, within the network itself, the size of the genomic bins dictates what path lengths correspond to long-range interactions. This influences the length of sampling motifs chosen for the MCMC sampling step—the sampled paths should be long enough to capture long-range interactions. Paths of length 21 nodes result in  $21 \times 500 = 10$ , 500 bases in the chromosome, which in turn amounts to a length of approximately 10, 000 bases. Additionally, the choice for the pathlength also controls the trade-off between the number of representatives and their size. With a choice of path-length as above, we have to draw 20, 000 MCMC samples to cover all the nodes

(chromatin fragments) in the dataset. This is evidenced by Fig D in <u>S1 Text</u> which plots the number of MCMC samples needed for given percentages of node coverage.

Similarly, the choice for K, the number of dictionary elements used, also depends on the dataset. Promoters and enhancers only constitute a very small fraction of the entire length of noncoding DNA. Studies indicate the existence of 10, 000 to 12, 000 such regions in the *Drosophila* genome, with each region being 100-1000 bases in length [34,35]. Working with the upper range of values, we arrive at a total length upper bounded by  $12,000 \times 1,000 = 12,000,000$  bases for the promoter/enhancer regions (one should compare this to the total length of the genome, which equals 180,000,000 bases). With K=25 for each of the 4 chromosomes, the dictionary elements will cover approximately  $4 \times 25 \times 10 \times 10,000 = 10,000,000$  bases which is close to the (loose) upper-bound estimate for the total length of the promoter/enhancer regions.

As a final remark, performing a multidimensional grid search for hyperparameters may be computationally prohibitive. Also, the procedure outlined above relies on solid biological side-information.

MCMC sampling for initialization, as well as for subsequent online optimization steps, was performed before running the online optimization process to improve the efficiency of our implementation. We sampled 20, 000 subnetworks from each of the four chromosomes to ensure sufficient coverage of the input network. From this pool of subnetworks, we randomly selected 500 subnetworks to initialize our dictionaries, ensuring that each dictionary element had at least 10 representatives (which suffice to get quality initializations for the dictionary elements themselves). Each online step involved sampling an additional subnetwork and we iterated this procedure up to 1 million times, as needed for convergence (see Fig 1A).

At this point, it is crucial to observe that the dictionary elements learned by online cvxNDL effectively capture *long-range interactions* because each dictionary element may include distal genomic regions that are not adjacent in the genomic order. In other words, the diagonal entries of our dictionary elements *do not exclusively represent consecutive genomic regions* as in standard chromatin contact maps; instead, they may include *both* nonconsecutive (long-range) and consecutive (short-range, adjacent) interactions. This point is explained in detail in Fig 4. Another relevant remark is that without the convexity constraint, dictionary element entries could not have been meaningfully mapped back (associated) to genomic regions and viewed as *real physical interactions between genomic loci*.

The dictionary elements generated from the *Drosophila* ChIA-Drop data for chr2L, chr2R, chr3L, and chr3R using the online cvxNDL method are shown in Fig 5. Each subplot corresponds to one chromosome and has 25 dictionary elements ordered with respect to their *importance scores*, capturing the relevance and frequency of use of the dictionary element, and formally defined in the Methods Section. Each element is color-coded based on the genomic locations covered by their representatives. Hence, dictionary elements represent combinations of experimentally observed interaction patterns, uniquely capturing the significance of the genomic locations involved in the corresponding interactions. We also report the density and median distance between all consecutive pairs of interacting loci (connected nodes) of all dictionary elements in Tables B and C in S1 Text.

Note that our algorithm is the first method for online learning of convex (interpretable) network dictionaries. We can therefore only compare its *representation accuracy* to that of nonnegative matrix factorization (NMF), convex matrix factorization (CMF), and online network dictionary learning (online NDL). A comparison of the dictionaries formed through online cvxNDL and the aforementioned methods for chr2L is provided in Fig 6.

Classical NMF does not allow the mapping of results back to real interacting genomic regions. While the dictionary elements obtained via CMF are interpretable, they tend to mostly

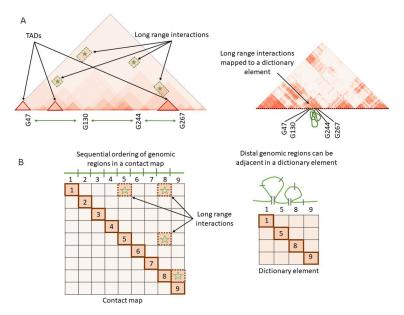


Fig 4. A dictionary element, represented as a matrix, consists of both proximal and distal interacting genomic regions. The elements on the diagonal are not necessarily indexed by adjacent (consecutive) genomic fragments, as explained by the example in the second row. There, off-diagonal long-range interactions in the  $9 \times 9$  matrix are included in a  $3 \times 3$  dictionary element whose diagonal elements are not in consecutive order.

comprise widely spread genomic regions since they do not use the network information. The dictionary elements generated by online cvxNDL have smaller yet relevant spreads that are more likely to capture meaningful long-range interactions. In contrast to online cvxNDL, both NMF and CMF are not scalable to large datasets, rendering them unsuitable for handling current and future high-resolution datasets such as those generated by ChIA-Drop. Compared to online NDL, online cvxNDL also has a more balanced distribution of importance scores. For example, in Fig 6B, dict\_0 has score 0.459, while the scores in Fig 6D are all  $\leq$  0.085. Moreover, akin to standard NMF, NDL fails to provide interpretable results since the dictionary elements cannot be mapped back to real interacting genomic loci.

Note that our approach is inherently an NMF-based method adapted for networks to ensure scalability, via its online nature, and interpretability, based on its convexity constraints. Besides scalability and interpretability, all the limitations of general NMF methods carry over to our method. For example, NMF approaches can be sensitive to initialization. Selection of the number of elements (or the rank of NMF) requires domain knowledge as well as heuristic search and testing. A wrong choice of the rank can lead to underfitting or overfitting the data. Furthermore, NMF does not guarantee a unique solution.

Results for other chromosomes are reported in Section D in S1 Text. Recall that both online cvxNDL and online NDL use a k-path as the template.

### **Reconstruction accuracy**

Once a dictionary is constructed, one can use the network reconstruction algorithm from [15] to recover a subnetwork or the whole network by locally approximating subnetworks via dictionary elements. The accuracy of approximation in this case measures the "expressibility" of the dictionary with respect to the network. All methods, excluding randomly generated dictionaries used for illustrative purposes only, can accurately reconstruct the input network. For a

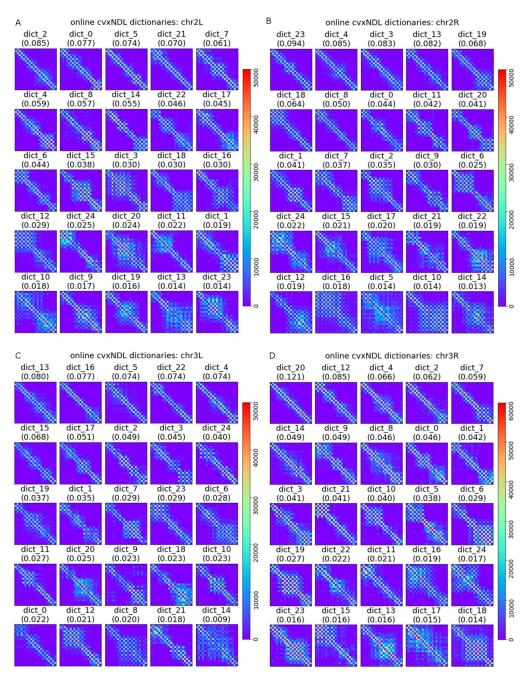


Fig 5. Dictionary elements for *Drosophila* chromosomes 2L, 2R, 3L and 3R obtained using online cvxNDL. Each subplot contains 25 dictionary elements for the corresponding chromosome and each block in the subplots corresponds to one dictionary element. The elements are ordered by their importance score. Note that the "diagonals" in the dictionary elements do not exclusively represent localized topologically associated domains (TADs) as in standard chromatin contact maps; instead, they can also capture long-range interactions. This is due to the fact that the indices of the dictionary element matrices represent genomic regions that may be far apart in the genome. In contrast, standard contact maps have indices that correspond to continuously ordered genomic regions, so that the diagonals truly represent TADs (see Fig 4). The color-code captures the actual locations of the genomic regions involved in the representatives and their dictionary elements. The most interesting dictionary elements are those that contain both dark blue, light blue/green, and red colors (since they involve long-range interactions). This is especially the case for chr3L and chr3R.

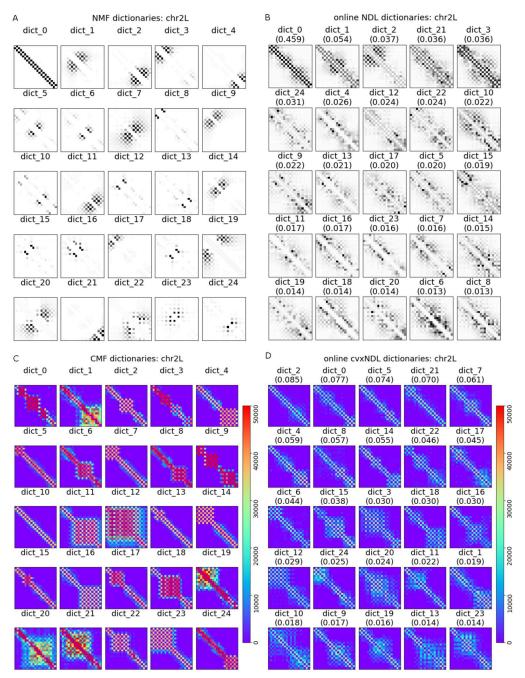


Fig 6. Dictionary elements for *Drosophila* chromosome chr2L generated by (A) NMF, (B) online NDL, (C) CMF and (D) online cvxNDL. NMF and CMF are learned off-line, using a total of 20, 000 samples. Note that these algorithms do not scale and cannot work with larger number of samples such as those used in online cvxNDL. The color-coding is performed in the same manner as for the accompanying online cvxNDL results. Columns of the dictionary elements in the second row are color-coded based on the genome locations of the representatives. As biologically meaningful locations can be determined only via convex methods, the top row corresponding to NMF and online NDL results is black-and-white.

**Table 1. Average Precision Recall for different DL methods, for all chromosomes as well as SBM synthetic datasets.** Methods that return interpretable dictionaries are indicated by the superscript *i* while methods that are scalable to large datasets are indicated by the superscript *s*. Online cvxNDL is both interpretable and scalable while maintaining performance on par with other noninterpretable and nonscalable methods.

	chr2L	chr2R	chr3L	chr3R	Synthetic (SBM)
Online cvxNDL <sup>i, s</sup>	0.9954	0.9986	0.9830	0.9876	0.9747
Online NDL <sup>s</sup>	0.9955	0.9986	0.9834	0.9880	0.9728
NMF	0.9952	0.9985	0.9829	0.9873	0.9774
$\overline{\mathrm{CMF}^i}$	0.9951	0.9985	0.9824	0.9870	0.9731
Random Dict.	0.0007	0.2547	0.5276	0.0796	0.1922

quantitative assessment, the average precision-recall score for all methods is plotted in Table 1. As expected, random dictionaries have the lowest scores across all chromosomes, while all other methods are of comparable quality. This means that interpretable methods, such as our online cvxNDL, do not introduce representation distortions (CMF also learns interpretable dictionaries; however, it is substantially more expensive computationally when compared to our method but does not ensure that network topology is respected). A zoomed-in sample-based reconstruction result for chr2L is shown in Fig H in S1 Text, while the reconstruction results for the entire contact maps of chr2L, chr2R, chr3L, and chr3R are available in Figs I-L in S1 Text. Additionally, for synthetic (Stochastic Block Model (SBM)) data, Fig 7 shows the reconstructed adjacency matrices for various dictionary learning methods, further confirming the validity of findings for the chromatin data. More detailed results for synthetic SBM data are available in Section C in S1 Text.

# Gene Ontology enrichment analysis

As each dictionary element is associated with a set of representatives that correspond to real observed subnetworks, their nodes can be mapped back to actual genomic loci. This allows one to create lists of genes covered by at least one node included in the representatives.

To gain insights into the functional annotations of the genes associated with the dictionary elements, we conducted a Gene Ontology (GO) enrichment analysis using the annotation category "Biological Process" from https://urldefense.com/v3/\_http://geneontology.org\_\_;!! DZ3fjg!4VWHhuROFHcJ1bWTZ8pNxUn75T-K3BfsdTvxM1iU1hXmSGX84JcRsXyIZZS0k 5Iaub9yNiansT9FS12EO52\_OaGhnYs\$, with the reference list *Drosophila Melanogaster*. This analysis was performed for each dictionary element. Our candidate set for enriched GO terms was selected with a false discovery rate (FDR) threshold of < 0.05. Note that the background genes used for comparison are all genes from all chromosomes (the default option). We also utilized the hierarchical structure of GO terms [36], where terms are represented as nodes in a directed acyclic graph, and their relationships are described via arcs in the digraph (i.e., each "child" GO term is more specific than its "parent" term and where one child may have multiple parents).

We further refined our results by running additional processing steps. For each GO term, we identified all the paths between the term and the root node and then removed any intermediate parent GO term from the enriched GO terms set. By iteratively performing this filtering process for each dictionary element, we created a list of the most specific GO terms associated with each element. More details about the procedure are available in Section F in S1 Text.

We report the most frequently enriched GO terms for each chromosome, along with the corresponding dictionary elements exhibiting enrichment for chr3R in Fig 8. The results for other chromosomes are available in Tables D, E, and F in S1 Text. Notably, the most frequent

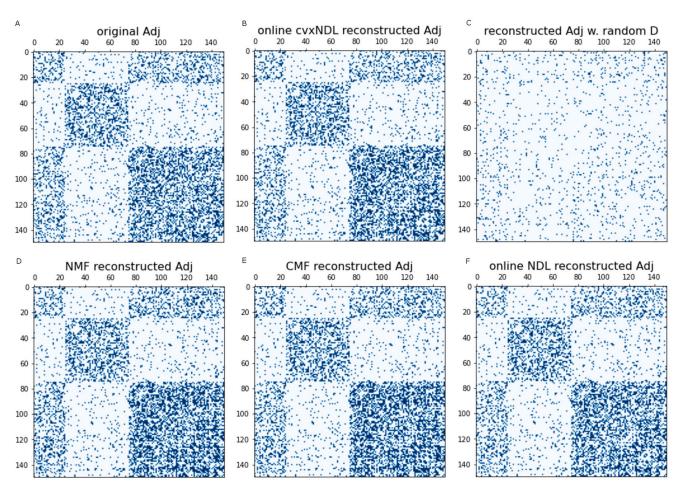


Fig 7. Original adjacency matrix and reconstructed adjacency matrices based on different DL methods, for an example Stochastic Block Model (SBM), including randomly selected dictionary elements. Both the x and y axes in the figures index the nodes of the synthetic network generated by the stochastic block model (SBM). The nodes are reorganized to highlight the underlying community structure. For a more quantitative analytical accuracy comparisons, see Table 1.

GO terms are related to regulatory functions, reflecting the significance of RNA Polymerase II. We also observe that dictionary elements for chr2L and chr2R are enriched in GO terms associated with reproduction and embryonic development. Similarly, chr3L and 3R are enriched in GO terms for blood circulation and responses to heat and cold.

We report the number of GO terms associated with each dictionary element, along with their importance scores in Tables J-M in S1 Text. Dictionary elements with higher importance scores tend to exhibit a larger number of enriched GO terms while dictionary elements with 0 enriched GO terms generally have small importance scores.

Using the entire genome as the reference is an accepted approach for GO analysis. However, it can introduce a bias due to differences in the chromosomal architectures of various chromosomes. We therefore performed an additional GO analysis where the genes within the pertinent chromosome, rather than the whole genome, are used as a reference. We implemented a Bonferroni correction and set the FDR to 0.05 (note that the results depend on the multiple-hypothesis testing correction method used). The total number of enriched GO terms across all online cvxNDL dictionaries for each of the 4 chromosomes 2L, 2R, 3L, and 3R equals 36, 19, 21, and 54, respectively.

Most frequent GO term		Top 3 dictionaries		
(GO:0001819) Positive regulation of cytokine production	7	dict. 20 dict. 7 dict. 9 (0.059) (0.049) $\rho$ =0.126, 0.146, 0.157 $d_{\rm med}$ =12791, 12830, 11930		
(GO:0008015) Blood circulation	7	dict, 20 dict, 12 dict, 4 (0.055) (0.056) (0.056) $\rho$ =0.126, 0.142, 0.138 $d_{med}$ =12791, 13455, 13674		
(GO:0045948) Positive regulation of translational initiation	5	dict. 20 dict. 4 dict. 14 (0.066) (0.049) $\rho$ = 0.126, 0.138, 0.162 $d_{med}$ = 12791, 13674, 12572		
(GO:0042177) Negative regulation of protein catabolic process	5	$\begin{array}{c} \text{dict, 20} \\ \text{(0.171)} \\ \text{(0.085)} \\ \end{array}$ $\begin{array}{c} \text{dict, 4} \\ \text{(0.085)} \\ \text{(0.086)} \\ \end{array}$ $\rho = 0.126, 0.142, 0.138 \ d_{\text{med}} = 12791, 13455, 13674 \end{array}$		
(GO:0043065) Positive regulation of apoptotic process	4	$ \begin{array}{c} \text{dict, 20} \\ \text{(0.171)} \\ \text{(0.059)} \\ \\ \rho = 0.126, 0.146, 0.179 \\ \end{array} \begin{array}{c} \text{dict, 7} \\ \text{(0.041)} \\ \\ $		

Fig 8. The 5 most enriched GO terms for genes covered by dictionary elements from chr3R. Column '#' indicates the number of dictionary elements that show enrichment for the given GO term. Also reported are up to 3 dictionary elements with the largest importance score in the dictionary, along with the "density"  $\rho$  of interactions in the dictionary element (defined in the Methods section) and median distance  $d_{\rm med}$  of all adjacent pairs of nodes in its representatives.

# **RNA-Seq coexpression analysis**

The ChIA-Drop dataset [17] used in our analysis was accompanied by a single noisy RNA-Seq replicate. To address this issue, we retrieved 20 collections of RNA-Seq data corresponding to untreated S2 cell lines of *Drosophila Melanogaster* from the Digital Expression Explorer (DEE2) repository. DEE2 provides uniformly processed RNA-Seq data sourced from the publicly available NCBI Sequence Read Archive (SRA) [23]. The list of sample IDs is available in Table N in S1 Text.

To ensure consistent normalization across all samples, we used the trimmed mean of M values (TMM) method [37], available through the edgeR package [38]. This is of crucial importance when jointly analyzing samples from multiple sources. We selected the most relevant genes by filtering the list of covered genes and retaining only those with more than 95% overlap with the gene promoter regions, as defined in the *Ensembl* genome browser. Subsequently, for each dictionary element, we collected all genes covered by it and then calculated the pairwise Pearson correlation coefficient of expressions of pairs of genes in the set. To visualize the underlying coexpression clusters within the genes, we performed hierarchical clustering, the results of which are shown in Section G in S1 Text and discussed next.

Additionally, we conducted control experiments by constructing dictionary elements through random sampling of genes from the list of all genes on each of the chromosomes. For these randomly constructed dictionaries, we carried out a coexpression analysis as described above. We observed that the mean of coexpressions of all pairs of genes in a randomly constructed dictionary element is significantly lower compared to the mean of the online cvxNDL dictionary elements. Specifically, for dictionary elements generated using online cvxNDL, the mean coexpression values for all pairs of genes covered by the 25 dictionary elements, and for each of the four chromosomes, 2L, 2R, 3L, and 3R, were found to be 0.419, 0.383, 0.411, and 0.407, respectively. The corresponding values for randomly constructed dictionaries were

found to be 0.333, 0.329, 0.323, and 0.337, respectively. To determine if these differences are statistically significant, we employed the two-sample Kolmogorov-Smirnov test [39], comparing the empirical cumulative distribution functions (ECDFs) of pairwise coexpression values of the learned and randomly constructed dictionaries. The null hypothesis used was "the two sets of dictionary elements are drawn from the same underlying distribution." The null hypotheses for all four chromosomes were rejected, with p-values equal to  $3.6 \times 10^{-9}$ ,  $8.5 \times 10^{-6}$ ,  $3.6 \times 10^{-9}$ , and  $2.5 \times 10^{-7}$  for chr2L, chr2R, chr3L, and chr3R, respectively (see Fig 9). This indicates that the learned dictionary elements indeed capture meaningful biological patterns of chromatin interactions.

To further evaluate our results, we also examined the well-documented R1-R4 and T1-T4 TAD interactions on chr2L, reported in [17]. The results of the coexpression analysis for these genomic regions are reported in Fig 10. The mean pairwise correlation between genes belonging to the R1-R4 genomic regions equals 0.422, which is comparable to the mean value 0.419 of the results obtained via online cvxNDL. We also calculated the intersection of the set of genes within the R1-R4 genomic regions and the set of genes covered by online cvxNDL dictionary elements identified for chr2L. We observed that the top 5 online cvxNDL dictionary elements cover 38 out of 85 genes in the R1-R4 genomic regions. This is to be contrasted with the results for random dictionary elements, which cover only 7 genes. Table 2 describes these and related findings in more detail.

We also mapped genes covered by our dictionary elements onto nodes of the STRING protein-protein interaction network [24]. These mappings allow us to determine the confidence of pairwise gene interactions. These, and related results based on FlyMine [40] data, a large gene expression repository for *Drosophila Melanogaster*, are available in Section G in S1 Text.

The rationale behind the STRING analysis is that gene fragments that are in physical contact are likely to be involved in the same pathway. This hypothesis, as well as the hypothesis we

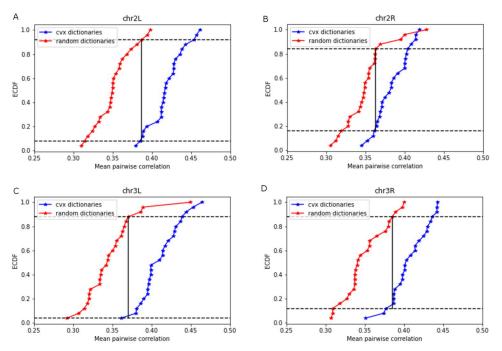


Fig 9. Empirical cumulative distribution functions (ECDF) of mean pairwise coexpressions of genes covered by random and online cvxNDL dictionary elements ((a) for chr2L, (b) for chr2R, (c) for chr3L and (d) for chr3R). The results are based on the two-sample Kolmogorov-Smirnov test, and the null hypothesis described in the main text.

https://doi.org/10.1371/journal.pcbi.1012095.g009

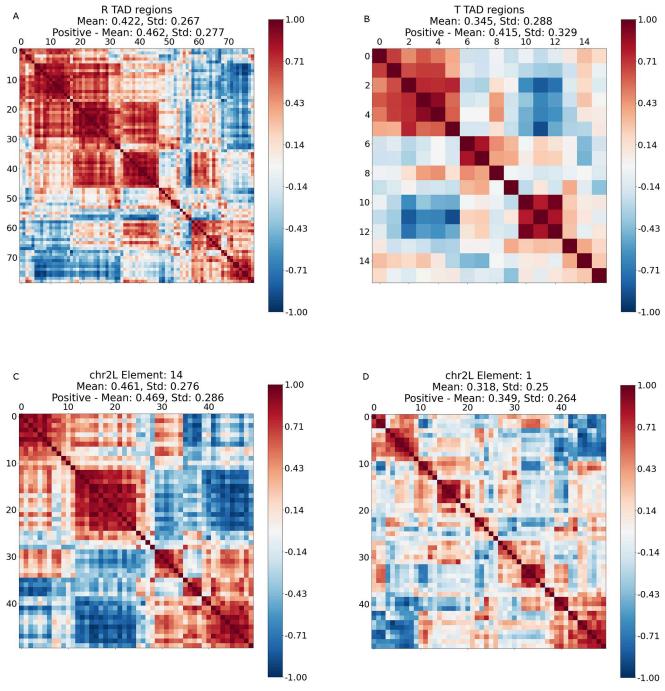


Fig 10. Pairwise coexpression of genes covered by (a) the R1-R4 genomic regions, (b) the T1-T4 genomic regions, (c) an online cvxNDL dictionary element, and (d) a randomly constructed dictionary element. We calculated the mean and standard deviation of absolute pairwise coexpression values, and the mean and standard deviation of coexpression values specifically for all positively correlated gene pairs. The mean coexpression values within TADs and dictionary elements are similar to each other and generally higher than those of randomly constructed dictionary elements. The x and y axis index genes that belong to the respective TAD regions or a specific dictionary element. Note that the plot (b) is of coarser resolution due to the small number of genes covered when compared to the cases (a), (c), (d).

Table 2. Intersection between the set of genes within the R1-R4 genomic regions and the sets of genes covered by online cvxNDL dictionary elements for chr2L. We determined the sizes of the intersections of the set of genes covered by each dictionary element and the genes in the R1-R4 genomic region and arranged them in decreasing order. The top 5 dictionary elements in this order cumulatively contain 38 out of the 85 genes within the R1-R4 genomic regions. This is in sharp contrast with randomly generated dictionary elements, where the top 5 elements with maximum intersection cover only 7 genes.

	Online cvxNDL			Random		
	Dictionary element id	Intersection	Cumulative	Dictionary element id	Intersection	Cumulative
1	1	15	15	20	3	3
2	11	12	24	0	1	4
3	12	12	30	1	1	5
4	7	11	35	21	1	6
5	21	10	38	17	1	7

used for RNA-Seq based validation that genes in high-frequency contact regions are coexpressed is, still being investigated. While there is evidence to suggest that the formation of loops causes coexpression of its gene constituents, the dynamic nature of chromatin folding and the potential rewiring of chromatin contacts before transcription may make the relationship more complex [18]. This is especially the case when some promoters act as enhancers during transcription of proteins [41].

We provide further comparisons of CMF and online cvxNDL methods, the only two interpretable methods, below. To ensure a fair comparison, we select the top 10 samples with the largest convex weights that correspond to each of the CMF dictionary elements.

The total number of enriched GO terms across all online cvxNDL dictionaries for each of the 4 chromosomes 2L, 2R, 3L, and 3R is 36, 19, 21, and 54, while the numbers for CMF are 17, 7, 2, and 26, respectively. Furthermore, although fine-grained GO term comparison is not possible for the two sets (since there are many specializations of the same higher-level term), we still see that important higher-level GO terms—such as protein folding, response to stimuli, and metabolic and developmental processes—are shared by the two lists.

The mean pairwise co-expressions of genes covered by all 25 online cvxNDL and CMF dictionary elements for each of the four chromosomes analyzed (and their standard deviation) are shown in Table 3. Also, both CMF and online cvxNDL dictionaries significantly outperform random dictionary elements. Similarly, the mean confidence values of interactions retrieved from the filtered STRING PPI network are reported in Table O in S1 Text. The RNA-Seq and PPI network analysis indicates that the interpretable dictionaries from online cvxNDL and CMF both perform similarly, while only online cvxNDL can scale to larger datasets.

We also reconstructed the R1-R4 genomic regions identified in [17] using CMF dictionary elements. We observe that the top 5 dictionary elements, with the 5-highest importance scores, capture 15 of the 85 genes present in the R1-R4 genomic regions. This is to be compared to the 38 genes covered by the top 5 online cvxNDL dictionary elements and the 7 genes covered by

Table 3. The mean pairwise co-expressions of genes (and their standard deviations) covered by all 25 online cvxNDL and CMF dictionary elements for each of the four chromosomes.

Chromosome	Online cvxNDL	CMF
chr2L	0.419 (0.022)	0.407 (0.025)
chr2R	0.383 (0.020)	0.386 (0.017)
chr3L	0.411 (0.025)	0.413 (0.031)
chr3R	0.407 (0.024)	0.404 (0.026)

https://doi.org/10.1371/journal.pcbi.1012095.t003

the top 5 random dictionary elements. All 25 CMF dictionary elements together cover 45 genes, and this number is comparable to that of 54 genes covered by all 25 online cvxNDL dictionary elements. It is also significantly larger than the 11 genes covered by 25 random dictionary elements.

### **Methods**

### **Notation**

Sets of consecutive integers are denoted by  $[l] = \{1, ..., l\}$ . The symbol  $\mathbb N$  is reserved for the natural numbers. Capital letters are reserved for matrices (bold font) and random variables (RVs) (regular font). Vectors are denoted by lower-case underlined letters. For a matrix of dimension  $d \times n$  over the reals,  $\mathbf{A} \in \mathbb{R}^{d \times n}$ ,  $\mathbf{A}[i, :]$  is used to denote the  $i^{\text{th}}$  row and  $\mathbf{A}[i, :]$  the  $i^{\text{th}}$  column of  $\mathbf{A}$ . The entry in row i, column j is denoted by  $\mathbf{A}[i, j]$ . Similarly,  $\underline{x}[l]$  is used to denote the  $l^{\text{th}}$  coordinate of a deterministic vector  $\underline{x} \in \mathbb{R}^d$ . Furthermore, we use the standard notation for the  $\ell_1$  and Frobenius norm of matrices,  $\|\mathbf{A}\|_1 = \sum_{i,j} |\mathbf{A}[i,j]|$  and  $\|\mathbf{A}\|_F^2 = \sum_{i,j} \mathbf{A}[i,j]^2$ , respectively.

A network  $\mathcal{G} = ([n], \mathbf{A})$  is an ordered pair of sets, the node set [n], and the set of edges represented by their adjacency matrix  $\mathbf{A}$ . Our underlying assumption is that the network is connected, which means that every node can be reached from every other node. Also,  $\mathbf{A}[i,j] = \mathbf{A}[j,i] \in \{0,1\}$ , indicating the presence or absence of an undirected edge between nodes i,j. In addition,  $\operatorname{Col}(\mathbf{A})$  stands for the set of columns of  $\mathbf{A}$ , while  $\operatorname{cvx}(\mathbf{A})$  stands for the convex hull of  $\operatorname{Col}(\mathbf{A})$ .

### Online DL

We first formulate the online DL problem. Assume that N input data samples are generated by a random process and organized in matrices  $(\mathbf{X}_t)_{t\in\mathbb{N}}\in\mathbb{R}^{d\times N}$  indexed by time t. For N=1,  $\mathbf{X}_t$  reduces to a column vector that encodes a d-dimensional signal. Given an online, sequentially observed data stream  $(\mathbf{X}_t)_{t\in\mathbb{N}}$ , the goal is to find a sequence of dictionary matrices  $(\mathbf{D}_t)_{t\in\mathbb{N}}, \mathbf{D}_t \in \mathbb{R}^{d\times K}$ , and codes  $(\mathbf{\Lambda}_t)_{t\in\mathbb{N}}, \mathbf{\Lambda}_t \in \mathbb{R}^{K\times N}$ , such that when  $t\to\infty$  almost surely we have

$$\|\mathbf{X}_{t} - \mathbf{D}_{t} \mathbf{\Lambda}_{t}\|_{F}^{2} \to \min_{\mathbf{D}, \mathbf{\Lambda}} \mathbb{E}_{\mathbf{X}} \|\mathbf{X} - \mathbf{D} \mathbf{\Lambda}\|_{F}^{2}.$$
 (1)

The expected loss in Eq 1 can be minimized by iteratively updating  $\Lambda_t$  and  $\mathbf{D}_t$  every time a new data sample  $\mathbf{X}_t$  is observed. The approximation error of  $\mathbf{D}$  for a single data sample  $\mathbf{X}$  is chosen as

$$l(\mathbf{X}, \mathbf{D}) = \min_{\mathbf{\Lambda} \in \mathbb{D}^{K \times N}} \|\mathbf{X} - \mathbf{D}\mathbf{\Lambda}\|_F^2 + \lambda \|\mathbf{\Lambda}\|_1.$$
 (2)

The second term represents a sparsity-enforcing regularizer. Furthermore, the empirical  $f_t$  and surrogate loss  $\hat{f}_t$  for **D** are defined as

$$f_t(\mathbf{D}) = (1 - w_t)f_{t-1}(\mathbf{D}) + w_t l(\mathbf{X}_t, \mathbf{D}), t \ge 1,$$
 (3)

$$\hat{f}_t(\mathbf{D}) = (1 - w_t)\hat{f}_{t-1}(\mathbf{D}) + w_t(\|\mathbf{X}_t - \mathbf{D}\mathbf{\Lambda}\|_F^2 + \lambda\|\mathbf{\Lambda}\|_1), \tag{4}$$

where the weight  $w_t$  determines the sensitivity of the algorithm to the newly observed data. The online DL algorithm first updates the code matrix  $\Lambda_t$  by solving Eq.(2) with  $l(\mathbf{X}_t, \mathbf{D}_{t-1})$ ,

then updates the dictionary matrix  $\mathbf{D}_t$  by minimizing (4) via

$$\mathbf{D}_{t} = \underset{\mathbf{D} \in \mathbb{R}^{d \times r}}{\operatorname{arg \, min}} \left( \operatorname{Tr}(\mathbf{D} \mathbf{A}_{t} \mathbf{D}^{T}) - 2 \operatorname{Tr}(\mathbf{D} \mathbf{B}_{t}) \right), \tag{5}$$

where  $\mathbf{A}_t = (1 - w_t)\mathbf{A}_{t-1} + w_t\mathbf{\Lambda}_t\mathbf{\Lambda}_t^T$  and  $\mathbf{B}_t = (1 - w_t)\mathbf{B}_{t-1} + w_t\mathbf{\Lambda}_t\mathbf{X}_t^T$  are the aggregated history of the input data and their codes, respectively. For simplicity, we set  $w_t = \frac{1}{t}$ .

To add convexity constraints, we introduce for each dictionary element a *representative set* (region)  $\hat{\mathbf{X}}_t^{(i)} \in \mathbb{R}^{d \times N_i}$ ,  $i \in [K]$ , where  $N_i$  is the size of the representative set for dictionary element  $\mathbf{D}_t[:,i]$ , and  $N = \sum_{i=1}^K N_i$ . The representative set for a dictionary element is a small subcollection of real data samples observed up to time t that best explain the dictionary element they are assigned to. The set of representatives is updated after observing a sample, the inclusion of which provides a better estimate of the dictionary element compared to the previous set. Since the representative set is bounded in size, if a new sample is included, an already existing sample has to be removed (see Fig 2B). Formally, the optimization objective is of the form

$$\min_{\mathbf{D} \in \text{cvx}(\hat{\mathbf{X}}), \hat{\mathbf{X}}} \hat{f}_t(\mathbf{D}) = \min_{\mathbf{D} \in \text{cvx}(\hat{\mathbf{X}}), \hat{\mathbf{X}}} \left( 1 - \frac{1}{t} \right) \hat{f}_{t-1}(\mathbf{D}) + \frac{1}{t} \left( \|\mathbf{X}_t - \mathbf{D}\boldsymbol{\Lambda}_t\|_F^2 + \lambda \|\boldsymbol{\Lambda}_t\|_1 \right). \tag{6}$$

# MCMC sampling of subnetworks (sample generation)

For NDL, it is natural to let the columns of  $\mathbf{X}_t$  be vectorized adjacency matrices of N subnetworks. Hence one needs to efficiently sample meaningful subnetworks from a (large) network. In image DL problems, samples can be generated directly from the image using adjacent rows and columns. However, such a sampling technique cannot be applied to arbitrary network data. Selecting nodes along with their one-hop neighbors at random may produce subnetworks of vastly different sizes and the results do not capture meaningful long-range interactions. It is also difficult to trim such subnetworks to uniform sizes. Furthermore, sampling a fixed number of nodes uniformly at random from sparse networks produces disconnected subnetworks with high probability and is not an acceptable approach either.

To address these problems, we consider "subnetwork sampling" introduced in [14, 15] where we fix a template network  $F = ([k], \mathbf{A}_F)$  of k nodes and seek subnetworks induced by k nodes in the input network  $\mathcal{G}$ , with the constraint that the subnetwork *contains* (but does not necessarily equals) the template F topology. Given an input network  $\mathcal{G} = ([n], \mathbf{A})$  and a template network  $F = ([k], \mathbf{A}_F)$ , we define a set of homomorphisms as a vector of the form

$$\operatorname{Hom}(F,\mathcal{G}) = \left\{ \underline{x} : [k] \to [n] \middle| \prod_{1 \le i,j \le k} \mathbf{A}[\underline{x}[i],\underline{x}[j]]^{\mathbf{A}_{F}[i,j]} = 1 \right\},\tag{7}$$

where we by default assume that  $0^0 = 1$ . For each homomorphism  $\underline{x} \in \operatorname{Hom}(F, \mathcal{G})$ , denote its induced adjacency matrix by  $\mathbf{A}_{\underline{x}}$ , where  $\mathbf{A}_{\underline{x}}[a,b] = \mathbf{A}[\underline{x}[a],\underline{x}[b]]$ ,  $1 \le a,b \le k$ . The adjacency matrix  $\mathbf{A}_{\underline{x}}$  represents one sample from the input network  $\mathcal{G}$ . An example homomorphism is shown in Fig 1D, where the input network  $\mathcal{G}$  contains n = 9 nodes and the template network F is a star network that contains k = 4 nodes. One proper homomorphism in this case is  $\underline{x}[a] = 9$ ,  $\underline{x}[b] = 6$ ,  $\underline{x}[c] = 4$ ,  $\underline{x}[d] = 7$ , which gives rise to an adjacency matrix  $\mathbf{A}_{\underline{x}}$  as depicted. A homomorphism can be sampled using the rejection sampling algorithm presented in Section B, Algorithm A in S1 Text. Our choice of template network, as already mentioned, is a k-path, i.e., a path joining k nodes. Paths are a simple and natural choice for networks with long average path lengths, such as chromatin interaction networks. It is also the same choice of template

used in standard NDL. As a final remark, we note that a k-path homomorphism leads to a sample of dimension  $d = k^2$ , as we will flatten its  $k \times k$  adjacency matrix into a single vector.

Although rejection sampling can be used repeatedly to generate several homomorphisms, it is highly inefficient. To efficiently generate a sequence of sample adjacency matrices  $\mathbf{A}_{\underline{x}_t}$  from  $\mathcal{G}$ , the MCMC sampling algorithm is used instead, while rejection sampling is only used to initialize the MCMC algorithm.

Next, for a homomorphism  $\underline{x}_t$ , let  $\mathcal{N}[\underline{x}_t[1]]$  ( $\mathcal{N}$  for short) denote the set of neighbors of  $\underline{x}_t[1]$ . We first choose a node  $v \in \mathcal{N}$  from the neighborhood of  $\underline{x}_t[1]$  uniformly at random, i.e. with probability  $P(v) = \frac{1}{|\mathcal{N}|}$ . We also calculate the probability of acceptance  $\beta$  for the selected node v. For a k-path template used in our approach, the value of  $\beta$  is given by

$$\beta = \min \left\{ \frac{\sum_{c \in [n]} A^{k-1}[v, c]}{\sum_{c \in [n]} A^{k-1}[\underline{x}_t[1], c]} \frac{\sum_{c \in [n]} A[\underline{x}_t[1], c]}{\sum_{c \in [n]} A[v, c]}, 1 \right\},$$
(8)

following the guidelines from [14, 15].

Next, we draw a value  $u \in [0, 1]$  uniformly at random. If  $u < \beta$ , we accept  $\underline{x}_{(t+1)}[1] = \nu$ , otherwise we reject  $\nu$  and reset  $\underline{x}_{(t+1)}[1] = \underline{x}_t[1]$ . We then perform a directed random walk from  $\underline{x}_{t+1}[1]$  of length equal to k-1 to obtain  $\underline{x}_{(t+1)}[2], \ldots, \underline{x}_{(t+1)}[k]$ . An illustration of the sampling procedure is shown in Fig 1E, while the detailed algorithm is presented in Section B, Algorithm B in S1 Text.

## Online convex NDL (online cvxNDL)

We start by initializing the dictionary  $\mathbf{D}_0$  and representative sets  $\{\hat{\mathbf{X}}_0^{(i)}\}$ ,  $i \in [K]$ , for each dictionary element. The algorithm for initialization is presented in Section B, Algorithm C in S1 Text. After initialization, we perform iterative optimization to generate  $\mathbf{D}_t$  and  $\{\hat{\mathbf{X}}_t^{(i)}\}$ ,  $i \in [K]$ , to reduce the loss at round t. At each iteration, we use MCMC sampling to obtain a k-node random subnetwork as sample  $\mathbf{X}_t$ , and then update the codes  $\Lambda_t$  based on the dictionary  $\mathbf{D}_{t-1}$  by solving the optimization problem in Eq.(2). Then we assign the current sample to a representative set of the closest dictionary element, say  $\mathbf{D}_{t-1}[:,j]$ , and jointly update its representative set  $\hat{\mathbf{X}}_t^{(j)}$  and all dictionaries  $\mathbf{D}_t$  as shown in Fig 2B. The iterative update algorithm for online cvxNDL is presented in Section B, Algorithm D in S1 Text.

The output of the algorithm is a dictionary matrix  $\mathbf{D}_T \in \mathbb{R}^{k^2 \times K}$ , where each column is a flattened vector of a dictionary element of size  $k \times k$ , and the representative sets  $\{\hat{\mathbf{X}}_T^{(i)}\}$ ,  $i \in [K]$ , for each dictionary element. Each representative set  $\hat{\mathbf{X}}_T^{(i)} \in \mathbb{R}^{k^2 \times N_i}$  contains  $N_i$  history-sampled subnetworks from the input network as its columns which are called the representatives of the dictionary element. The convex hull of all representatives of a dictionary element forms the representative region of the dictionary element. We can easily convert both the dictionary elements and representatives back to  $k \times k$  adjacency matrices. Due to the added convexity constraint, each dictionary element  $\mathbf{D}_T[:,j]$  at the final step T has the *interpretable* form:

$$\mathbf{D}_{T}[:,j] = \sum_{i \in [N_{j}]} w_{j,i} \hat{\mathbf{X}}_{T}^{(j)}[:,i], \quad \text{s.t. } \sum_{i \in [N_{j}]} w_{j,i} = 1, w_{j,i} \ge 0, i \in [N_{j}], j \in [K].$$
(9)

The weight  $w_{j,i}$ ,  $i \in [N_j]$ , is the *convex coefficient* of the  $i^{th}$  representative of dictionary element  $\mathbf{D}_T[:,j]$ . Dictionary elements learned from the data stream can be used to reconstruct the input network by multiplying it with the dictionary element weights from  $\mathbf{Eq}$  (2). The  $j^{th}$  index of the weight vector corresponds to the contribution of dictionary element  $\mathbf{D}_{T-1}[:,j]$  to the reconstruction. Similarly to what was done in [15], we can also define the *importance score* for each

dictionary element as

$$\gamma(i) = \frac{\mathbf{A}_t[i,i]^2}{\sum_{j \in [K]} \mathbf{A}_t[j,j]^2}.$$
 (10)

We use the importance scores, as described in the previous sections, to determine the most frequently used interactions in the dictionary construction, as well as the most typical and important long-range interactions.

To conclude, we point out that the *density*  $\rho$  of interactions in a dictionary element is defined as

$$\rho = \frac{1}{k^2} \sum_{i,j=1}^k \mathbf{D}_T[i,j].$$

# Supporting information

S1 Text. Supplement PDF. Supplemental material, including figures and tables, is available in the Supplement file. The online cvxNDL code and test datasets are available at: <a href="https://urldefense.com/v3/\_https://github.com/rana95vishal/chromatin\_DL/\_;!!DZ3fjg!4VWHhuROFHcJ1bWTZ8pNxUn75T-K3BfsdTvxM1iU1hXmSGX84JcRsXyIZZS0k5Iaub9yNiansT9FS12EO52\_XsbpA\_s\$. A tool that enables readers with color-blindness to view the images using a more appropriate color palette is described at the end of the Supplement. (PDF)

# **Acknowledgments**

The authors gratefully acknowledge many useful discussions with Dr. Yijun Ruan.

### **Author Contributions**

Conceptualization: Minji Kim, Olgica Milenkovic.

**Data curation:** Vishal Rana, Jianhao Peng, Chao Pan, Albert Cheng, Minji Kim, Olgica Milenkovic.

**Formal analysis:** Vishal Rana, Jianhao Peng, Chao Pan, Hanbaek Lyu, Minji Kim, Olgica Milenkovic.

Funding acquisition: Olgica Milenkovic.

Investigation: Vishal Rana, Jianhao Peng, Chao Pan, Hanbaek Lyu, Minji Kim, Olgica Milenkovic.

**Methodology:** Vishal Rana, Jianhao Peng, Chao Pan, Hanbaek Lyu, Minji Kim, Olgica Milenkovic.

Project administration: Olgica Milenkovic.

Resources: Olgica Milenkovic.

Software: Vishal Rana, Jianhao Peng, Chao Pan, Minji Kim.

Supervision: Minji Kim, Olgica Milenkovic.

Validation: Vishal Rana, Jianhao Peng, Chao Pan, Minji Kim, Olgica Milenkovic.

Visualization: Vishal Rana, Jianhao Peng, Chao Pan, Minji Kim, Olgica Milenkovic.

Writing – original draft: Vishal Rana, Jianhao Peng, Chao Pan, Minji Kim, Olgica Milenkovic.

Writing – review & editing: Vishal Rana, Jianhao Peng, Chao Pan, Hanbaek Lyu, Albert Cheng, Minji Kim, Olgica Milenkovic.

### References

- Elad M, Aharon M. Image denoising via sparse and redundant representations over learned dictionaries. IEEE Transactions on Image processing. 2006; 15(12):3736–3745. https://doi.org/10.1109/TIP. 2006.881969 PMID: 17153947
- Mairal J, Elad M, Sapiro G. Sparse representation for color image restoration. IEEE Transactions on image processing. 2007; 17(1):53–69. https://doi.org/10.1109/TIP.2007.911828
- Cichocki A, Lee H, Kim YD, Choi S. Non-negative matrix factorization with α-divergence. Pattern Recognition Letters. 2008; 29(9):1433–1440. https://doi.org/10.1016/j.patrec.2008.02.016
- Ye M, Qian Y, Zhou J. Multitask sparse nonnegative matrix factorization for joint spectral–spatial hyperspectral imagery denoising. IEEE Transactions on Geoscience and Remote Sensing. 2014; 53 (5):2621–2639. https://doi.org/10.1109/TGRS.2014.2363101
- Lu H, Sang X, Zhao Q, Lu J. Community detection algorithm based on nonnegative matrix factorization and pairwise constraints. Physica A: Statistical Mechanics and its Applications. 2020; 545:123491. https://doi.org/10.1016/j.physa.2019.123491
- Zhu X, Ching T, Pan X, Weissman SM, Garmire L. Detecting heterogeneity in single-cell RNA-Seq data by non-negative matrix factorization. PeerJ. 2017; 5:e2888. https://doi.org/10.7717/peerj.2888 PMID: 28133571
- Shao C, Höfer T. Robust classification of single-cell transcriptome data by nonnegative matrix factorization. Bioinformatics. 2017; 33(2):235–242. <a href="https://doi.org/10.1093/bioinformatics/btw607">https://doi.org/10.1093/bioinformatics/btw607</a> PMID: 27663498
- Zhang S, Chasman D, Knaack S, Roy S. In silico prediction of high-resolution Hi-C interaction matrices. Nature communications. 2019; 10(1):1–18. <a href="https://doi.org/10.1038/s41467-019-13423-8">https://doi.org/10.1038/s41467-019-13423-8</a> PMID: 31811132
- Paatero P, Tapper U. Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. Environmetrics. 1994; 5(2):111–126. https://doi.org/10.1002/env. 3170050203
- Paatero P. Least squares formulation of robust non-negative factor analysis. Chemometrics and intelligent laboratory systems. 1997; 37(1):23–35. https://doi.org/10.1016/S0169-7439(96)00044-5
- Ding CH, Li T, Jordan MI. Convex and semi-nonnegative matrix factorizations. IEEE transactions on pattern analysis and machine intelligence. 2010; 32(1):45–55. <a href="https://doi.org/10.1109/TPAMI.2008.277">https://doi.org/10.1109/TPAMI.2008.277</a>
   PMID: 19926898
- 12. Mairal J, Bach F, Ponce J, Sapiro G. Online learning for matrix factorization and sparse coding. Journal of Machine Learning Research. 2010; 11(Jan):19–60.
- **13.** Peng J, Milenkovic O, Agarwal A. Online convex matrix factorization with representative regions. In: Advances in Neural Information Processing Systems; 2019. p. 13242–13252.
- Lyu H, Memoli F, Sivakoff D. Sampling random graph homomorphisms and applications to network data analysis. Journal of machine learning research. 2023; 24(9):1–79.
- Lyu H, Needell D, Balzano L. Online matrix factorization for Markovian data and applications to Network Dictionary Learning. Journal of Machine Learning Research. 2020; 21(251):1–49.
- Lyu H, Kureh YH, Vendrow J, Porter MA. Learning low-rank latent mesoscale structures in networks. Nature Communications. 2024; 15(1):224. <a href="https://doi.org/10.1038/s41467-023-42859-2">https://doi.org/10.1038/s41467-023-42859-2</a> PMID: 38172092
- Zheng M, Tian SZ, Capurso D, Kim M, Maurya R, Lee B, et al. Multiplex chromatin interactions with single-molecule precision. Nature. 2019; 566(7745):558–562. https://doi.org/10.1038/s41586-019-0949-1 PMID: 30778195
- Li G, Ruan X, Auerbach RK, Sandhu KS, Zheng M, Wang P, et al. Extensive promoter-centered chromatin interactions provide a topological basis for transcription regulation. Cell. 2012; 148(1-2):84–98. https://doi.org/10.1016/j.cell.2011.12.014 PMID: 22265404

- Tang Z, Luo OJ, Li X, Zheng M, Zhu JJ, Szalaj P, et al. CTCF-mediated human 3D genome architecture reveals chromatin topology for transcription. Cell. 2015; 163(7):1611–1627. <a href="https://doi.org/10.1016/j.cell.2015.11.024">https://doi.org/10.1016/j.cell.2015.11.024</a> PMID: 26686651
- Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. science. 2009; 326(5950):289–293. https://doi.org/10.1126/science.1181369 PMID: 19815776
- Li G, Fullwood MJ, Xu H, Mulawadi FH, Velkov S, Vega V, et al. ChIA-PET tool for comprehensive chromatin interaction analysis with paired-end tag sequencing. Genome biology. 2010; 11(2):R22. <a href="https://doi.org/10.1186/gb-2010-11-2-r22">https://doi.org/10.1186/gb-2010-11-2-r22</a> PMID: 20181287
- 22. Fullwood MJ, Liu MH, Pan YF, Liu J, Xu H, Mohamed YB, et al. An oestrogen-receptor-α-bound human chromatin interactome. Nature. 2009; 462(7269):58–64. https://doi.org/10.1038/nature08497 PMID: 19890323
- 23. Ziemann M, Kaspi A, El-Osta A. Digital expression explorer 2: a repository of uniformly processed RNA sequencing data. Gigascience. 2019; 8(4):giz022. https://doi.org/10.1093/gigascience/giz022 PMID: 30942868
- Szklarczyk D, Gable AL, Lyon D, Junge A, Wyder S, Huerta-Cepas J, et al. STRING v11: protein–protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. Nucleic acids research. 2019; 47(D1):D607–D613. https://doi.org/10.1093/nar/gky1131 PMID: 30476243
- Wang S, Zhang Q, He Y, Cui Z, Guo Z, Han K, et al. DLoopCaller: A deep learning approach for predicting genome-wide chromatin loops by integrating accessible chromatin landscapes. PLoS Computational Biology. 2022; 18(10):e1010572. https://doi.org/10.1371/journal.pcbi.1010572 PMID: 36206320
- Xie WJ, Qi Y, Zhang B. Characterizing chromatin folding coordinate and landscape with deep learning. PLoS computational biology. 2020; 16(9):e1008262. https://doi.org/10.1371/journal.pcbi.1008262 PMID: 32986691
- Zhang P, Wu Y, Zhou H, Zhou B, Zhang H, Wu H. CLNN-loop: a deep learning model to predict CTCF-mediated chromatin loops in the different cell lines and CTCF-binding sites (CBS) pair types. Bioinformatics. 2022; 38(19):4497–4504. https://doi.org/10.1093/bioinformatics/btac575 PMID: 35997565
- 28. Tian SZ, Li G, Ning D, Jing K, Xu Y, Yang Y, et al. MClBox: a toolkit for single-molecule multi-way chromatin interaction visualization and micro-domains identification. Briefings in Bioinformatics. 2022; 23 (6):bbac380. https://doi.org/10.1093/bib/bbac380 PMID: 36094071
- 29. Agarwal S, Lim J, Zelnik-Manor L, Perona P, Kriegman D, Belongie S. Beyond pairwise clustering. In: 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR). vol. 2. IEEE; 2005. p. 838–845.
- **30.** Zhou D, Huang J, Schölkopf B. Learning with hypergraphs: Clustering, classification, and embedding. Advances in neural information processing systems. 2006; 19.
- Li P, Milenkovic O. Inhomogeneous hypergraph clustering with applications. Advances in neural information processing systems. 2017; 30.
- 32. Kim M, Zheng M, Tian SZ, Lee B, Chuang JH, Ruan Y. MIA-Sig: multiplex chromatin interaction analysis by signal processing and statistical algorithms. Genome biology. 2019; 20(1):1–13. https://doi.org/10. 1186/s13059-019-1868-z PMID: 31767038
- Dekker J, Heard E. Structural and functional diversity of topologically associating domains. FEBS letters. 2015; 589(20):2877–2884. https://doi.org/10.1016/j.febslet.2015.08.044 PMID: 26348399
- 34. Hoskins RA, Landolin JM, Brown JB, Sandler JE, Takahashi H, Lassmann T, et al. Genome-wide analysis of promoter architecture in Drosophila melanogaster. Genome research. 2011; 21(2):182–192. https://doi.org/10.1101/gr.112466.110 PMID: 21177961
- FitzGerald PC, Sturgill D, Shyakhtenko A, Oliver B, Vinson C. Comparative genomics of Drosophila and human core promoters. Genome biology. 2006; 7:1–22. https://doi.org/10.1186/gb-2006-7-7-r53 PMID: 16827941
- Musen MA. The protégé project: a look back and a look forward. Al matters. 2015; 1(4):4–12. <a href="https://doi.org/10.1145/2757001.2757003">https://doi.org/10.1145/2757001.2757003</a> PMID: 27239556
- Robinson MD, Oshlack A. A scaling normalization method for differential expression analysis of RNAseq data. Genome biology. 2010; 11(3):1–9. <a href="https://doi.org/10.1186/gb-2010-11-3-r25">https://doi.org/10.1186/gb-2010-11-3-r25</a> PMID: 20196867
- Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. bioinformatics. 2010; 26(1):139–140. https://doi.org/10.1093/ bioinformatics/btp616 PMID: 19910308

- 39. Massey FJ Jr. The Kolmogorov-Smirnov test for goodness of fit. Journal of the American statistical Association. 1951; 46(253):68–78. https://doi.org/10.1080/01621459.1951.10500769
- Lyne R, Smith R, Rutherford K, Wakeling M, Varley A, Guillier F, et al. FlyMine: an integrated database for Drosophila and Anopheles genomics. Genome biology. 2007; 8(7):1–16. https://doi.org/10.1186/gb-2007-8-7-r129 PMID: 17615057
- **41.** Dao LT, Spicuglia S. Transcriptional regulation by promoters with enhancer function. Transcription. 2018; 9(5):307–314. https://doi.org/10.1080/21541264.2018.1486150 PMID: 29889606