

TFEA.ChIP: A tool kit for transcription factor binding site enrichment analysis capitalizing on ChIP-seq datasets

Laura Puente-Santamaria¹, Luis del Peso^{1,2,3*}

1 Departamento de Bioquímica, Universidad Autónoma de Madrid (UAM) and Instituto de Investigaciones Biomédicas 'Alberto Sols' (CSIC-UAM), 28029 Madrid, Spain

2 IdiPaz, Instituto de Investigación Sanitaria del Hospital Universitario La Paz, 28029 Madrid, Spain.

3 CIBER de Enfermedades Respiratorias (CIBERES), Instituto de Salud Carlos III, 28029 Madrid, Spain.

* luis.peso@uam.es

Abstract

The identification of transcription factor (TF) responsible for the co-regulation of a specific set of genes is a common problem in transcriptomics. With the development of TFEA.ChIP we aim to provide a tool to estimate and visualize TF enrichment in a set of differentially expressed genes that takes into account the wide variation in TF's behavior across different cell types and stimuli. To that end, ChIP-Seq experiments from the ENCODE Consortium and GEO Datasets were gathered, and a database linking TFs with the genes they interact with in each ChIP-Seq experiment was generated. In its current state, TFEA.ChIP covers 333 different transcription factors in 1122 ChIP-Seq experiments, with over 150 cell types being represented.

Include
parison
sum an

Author summary

[I believe this is required for research articles, but not software papers.]

Method
and Sof
cles req
submis
quiries.

Introduction

In the most simple scenario, the comparison of the transcriptome of cells or organisms in two conditions leads to the identification of a set of differentially expressed (DE) genes, and the underlying assumption is that one or a few TFs regulate the expression of those genes. Traditionally, the identification of relevant TFs has relied on the use of position weight matrices (PWMs) to predict transcription factor binding sites (TFBSs) proximal to the DE genes [1]. The comparison of predicted TFBS in DE versus a set of control genes, reveals factors that are significantly enriched in the DE gene set. The prediction of TFBS using these approaches have been useful to narrow down potential binding sites, but can suffer from high rates of false positives. In addition, this approach is limited by design to sequence-specific transcription factors (TF) and thus unable to identify cofactors that bind indirectly to target genes. To overcome these limitations we developed the R package TFEA.ChIP, which exploits the vast amount of publicly available ChIP-Seq datasets to determine TFBS proximal to a given set of genes and computes enrichment analysis based on this experimentally-derived rich information. Specifically TFEA.ChIP, uses information derived from the hundreds of ChIP-Seq

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

experiments from the ENCODE Consortium [2] expanded to include additional datasets contributed to GEO database [3] [4] by individual laboratories representing the binding sites of factors not assayed by ENCODE. The package includes a set of tools to customize the ChIP data, perform enrichment analysis and visualize the results. Herein we describe the main characteristics of the package and compare the results produced by TFEA.ChIP vs those generated by Opposum, an state of the art TFBS identification software based on PWMs [5]. Our data indicate that the results of TFEA.ChIP and Opposum are coincident for those datasets where Opposum identifies clear TFBS candidate(s). In addition, TFEA.ChIP identified enriched factors for some data sets where Opposum was unable to find a significant match.

Design and implementation

Database and algorithm

TFEA.ChIP package includes analysis and visualization tools intended for the identification of TFBS enriched in a set of DE genes. To this end the package uses experimental information derived from 1122 ChIP-seq datasets, generated by the ENCODE consortium and individual researchers, testing a total of 333 different human transcription factors in a variety of cell types and experimental conditions. Thus, this compiled database covers 20-24% of the 1,391 [6] to 1600 [7] transcription factors estimated to be encoded by the human genome and includes proteins from all the major classes of DNA binding domains (Fig 1).

The supplementary table S1 Table. contains the complete list of the datasets included in the package along with their GEO accession numbers. Although the package is mainly focused towards analyzing expression data generated from human cells, TFEA.ChIP includes the option to use datasets coming from experiments in mice, translating mouse gene names to their equivalent ID on the human genome. For the analysis, either the actual set of DE genes or a list of genes sorted according to their expression in the conditions under study, must be provided as an input.

- Analysis of the association of TFBS and differential expression from 2x2 tables recording the presence of binding sites for a given TF in DE and control genes. The statistical significance of the association for each factor determined by a Fisher's exact test.
- GSEA analysis, based on the core function of the GSEA algorithm for R [8] [9], GSEA.EnrichmentScore.

Nulla mi mi, Figvenenatis sed ipsum varius, volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, S1 Video vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

Results

Nulla mi mi, venenatis sed ipsum varius, Table 1 volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id

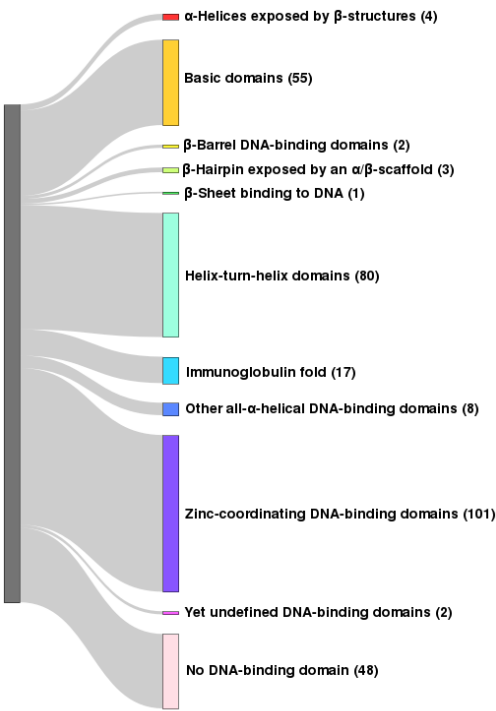


Fig 1. Structural diversity according to DNA-binding domains of the transcription factors included in the TFBS database. The 333 TFs included in TFEA.ChIP database were classified into families according to their DNA-binding domain composition. InterPro parent–child relationships between DNA-binding domains were used as the basis for TF family definition (Supplementary information S1 (PDF)). TFs with multiple DNA-binding domains were classified in each of their respective families. Families with less than five members were classified as ‘other’.

massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

Table 1. Table caption Nulla mi mi, venenatis sed ipsum varius, volutpat euismod diam.

Heading1				Heading2			
cell1row1	cell2 row 1	cell3 row 1	cell4 row 1	cell5 row 1	cell6 row 1	cell7 row 1	cell8 row 1
cell1row2	cell2 row 2	cell3 row 2	cell4 row 2	cell5 row 2	cell6 row 2	cell7 row 2	cell8 row 2
cell1row3	cell2 row 3	cell3 row 3	cell4 row 3	cell5 row 3	cell6 row 3	cell7 row 3	cell8 row 3

Table notes Phasellus venenatis, tortor nec vestibulum mattis, massa tortor interdum felis, nec pellentesque metus tortor nec nisl. Ut ornare mauris tellus, vel dapibus arcu suscipit sed.

LOREM and IPSUM nunc blandit a tortor

3rd level heading

Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur fringilla pulvinar lectus consectetur pellentesque. Quisque augue sem, tincidunt sit amet feugiat eget, ullamcorper sed velit. Sed non aliquet felis. Lorem ipsum dolor sit amet, consectetur

adipiscing elit. Mauris commodo justo ac dui pretium imperdiet. Sed suscipit iaculis mi at feugiat.

1. react
2. diffuse free particles
3. increment time by dt and go to 1

Sed ac quam id nisi malesuada congue

Nulla mi mi, venenatis sed ipsum varius, volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

- First bulleted item.
- Second bulleted item.
- Third bulleted item.

Discussion

Nulla mi mi, venenatis sed ipsum varius, Table 1 volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero [?].

Conclusion

CO₂ Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur fringilla pulvinar lectus consectetur pellentesque. Quisque augue sem, tincidunt sit amet feugiat eget, ullamcorper sed velit.

Sed non aliquet felis. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Mauris commodo justo ac dui pretium imperdiet. Sed suscipit iaculis mi at feugiat. Ut neque ipsum, luctus id lacus ut, laoreet scelerisque urna. Phasellus venenatis, tortor nec vestibulum mattis, massa tortor interdum felis, nec pellentesque metus tortor nec nisl. Ut ornare mauris tellus, vel dapibus arcu suscipit sed. Nam condimentum sem eget mollis euismod. Nullam dui urna, gravida venenatis dui et, tincidunt sodales ex. Nunc est dui, sodales sed mauris nec, auctor sagittis leo. Aliquam tincidunt, ex in facilisis elementum, libero lectus luctus est, non vulputate nisl augue at dolor. For more information, see S1 Appendix.

Supporting information

S1 Fig. Bold the title sentence. Add descriptive text after the title of the item (optional).

S2 Fig. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. 110
Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. 111
Curabitur fringilla pulvinar lectus consectetur pellentesque. 112

S1 File. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. 113
Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. 114
Curabitur fringilla pulvinar lectus consectetur pellentesque. 115

S1 Video. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. 116
Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. 117
Curabitur fringilla pulvinar lectus consectetur pellentesque. 118

S1 Appendix. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices 119
gravida. Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec 120
euismod ligula. Curabitur fringilla pulvinar lectus consectetur pellentesque. 121

S1 Table. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. 122
Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. 123
Curabitur fringilla pulvinar lectus consectetur pellentesque. 124

Acknowledgments 125

Cras egestas velit mauris, eu mollis turpis pellentesque sit amet. Interdum et malesuada 126
fames ac ante ipsum primis in faucibus. Nam id pretium nisi. Sed ac quam id nisi 127
malesuada congue. Sed interdum aliquet augue, at pellentesque quam rhoncus vitae. 128

References

1. Wasserman W, Sandelin A. Applied bioinformatics for the identification of regulatory elements. *Nature Reviews Genetics*. 2004;5:276. doi:10.1038/nrg1315.
2. Consortium EP. An integrated encyclopedia of DNA elements in the human genome. *Nature*. 2004;489:57–74. doi:10.1038/nature11247.
3. Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Research*. 2002;30:207–210.
4. Barrett T, et al. NCBI GEO: archive for functional genomics data sets - update. *Nucleic Acids Research*. 2013;41(D1):D991–D995. doi:10.1093/nar/gks1193.
5. Kwon AT, Arenillas DJ, Worsley Hunt R, Wasserman WW. oPOSSUM-3: advanced analysis of regulatory motif over-representation across genes or ChIP-Seq datasets. *G3 (Bethesda, Md)*. 2012;2(9):987–1002. doi:10.1534/g3.112.003202.
6. Vaquerizas JM, Kummerfeld SK, Teichmann SA, Luscombe NM. A census of human transcription factors: function, expression and evolution. *Nature Reviews Genetics*. 2009;10:252–263. doi:10.1038/nrg2538.
7. Lambert SA, Jolma A, Campitelli LF, Das PK, Yin Y, Albu M, et al. The Human Transcription Factors. *Cell*. 2018;172(4):650–665. doi:10.1016/j.cell.2018.01.029.

8. Subramanian A, Tamayo P, et al. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *PNAS*. 2005;102(43):15545–15550. doi:10.1073/pnas.0506580102.
9. Mootha VK, et al. PGC-1alpha-responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes. *Nature Genetics*. 2003;34:267–273. doi:10.1038/ng1180.