



## Combining transcription factor binding affinities with an open chromatin prior for accurate gene expression prediction

Florian Schmidt<sup>1,2</sup>, Nina Gasparoni<sup>3</sup>, Gilles Gasparoni<sup>3</sup>, Kathrin Gianmoena<sup>4</sup>, Cristina Cadenas<sup>4</sup>, Julia K. Polansky<sup>5</sup>, Peter Ebert<sup>2,6</sup>, Karl Nordström<sup>3</sup>, Matthias Barann<sup>7</sup>, Anupam Sinha<sup>7</sup>, Sebastian Fröhler<sup>8</sup>, Jieyi Xiong<sup>8</sup>, Azim Dehghani Amirabad<sup>1,2,6</sup>, Fatemeh Behjati Ardakani<sup>1,2</sup>, Barbara Hutter<sup>9</sup>, Gideon Zipprich<sup>10</sup>, Bärbel Felder<sup>10</sup>, Jürgen Eils<sup>10</sup>, Benedikt Brors<sup>9</sup>, Wei Chen<sup>8</sup>, Jan G. Hengstler<sup>4</sup>, Alf Hamann<sup>6</sup>, Thomas Lengauer<sup>2</sup>, Philip Rosenstiel<sup>7</sup>, Jörn Walter<sup>3</sup>, and Marcel H. Schulz<sup>1,2</sup>

*1: Cluster of Excellence for Multimodal Computing and Interaction, Saarland University, Saarbrücken, 66123, Germany.*

*2: Computational Biology & Applied Algorithmics, Max Planck Institute for Informatics, Saarbrücken, 66123, Germany.*

*3: Department of Genetics, University of Saarland, Saarbrücken, Germany.*

*4: Leibniz Research Centre for Working Environment and Human Factors IfADo, Ardeystrasse 67, Dortmund, 44139, Germany.*

*5: Experimental Rheumatology, German Rheumatism Research Centre, Berlin, Germany*

*6: International Max Planck Research School for Computer Science, Saarbrücken, 66123, Germany.*

*7: Institute of Clinical Molecular Biology, Christian-Albrechts-University, Kiel, Germany.*

*8: Berlin Institute for Medical Systems Biology, Max-Delbrück Center for Molecular Medicine, Berlin, Germany.*

*9: Applied Bioinformatics, Deutsches Krebsforschungszentrum, Heidelberg, Germany*

*10: Data Management and Genomics IT, Deutsches Krebsforschungszentrum, Heidelberg, Germany*

### Abstract

The binding and contribution of Transcription Factors (TFs) to cell specific gene expression is often deduced from open-chromatin measurements to avoid cost and labour intensive TF ChIP-seq assays. It is important to develop reliable and fast computational methods for accurate TF binding prediction in Open chromatin regions (ORCs). Here, we report a novel segmentation-based method, TEPIC, to predict TF binding by combining sets of OCRs with position weight matrices. TEPIC can be applied to various open chromatin data, e.g. DNase-Seq and NOMe-Seq, using either peaks or footprints as input data. TEPIC computes TF affinities as a quantitative measure (parameter) of TF binding strength and we show that low affinity binding sites predicted in this way improve performance over a simple presence/absence classification. Further, we show that while footprints called from OCRs capture most essential TF binding events, OCR peaks deliver the best prediction performance. Using machine learning techniques, we assessed the importance of individual TFs for gene expression and found that TEPIC scores nearly reach the quality of TF ChIP-seq data. Finally we show that TEPIC predicts all major known key transcriptional regulators in primary human hepatocytes and CD4+ T-cells emphasizing the reliability and applicability of our method.