# Data S4

# DNA shape features improve transcription factor binding site predictions *in vivo*
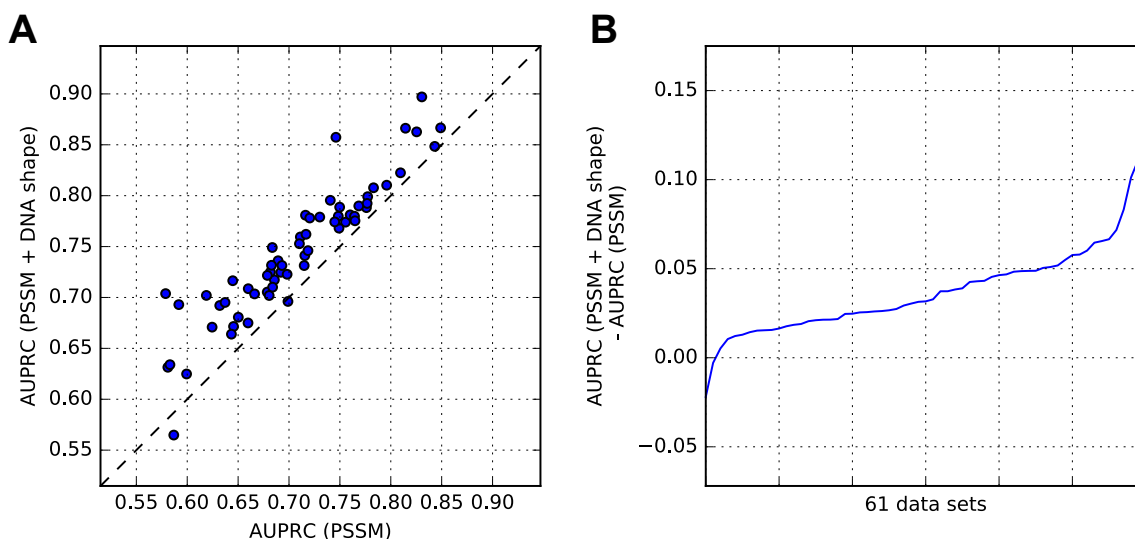
ANTHONY MATHELIER[1,2], BEIBEI XIN[3], TSU-PEI CHIU[3],
LIN YANG[3], REMO ROHS[3,*], AND WYETH W. WASSERMAN[1,*]

[1] Centre for Molecular Medicine at the Child and Family Research Institute,
Department of Medical Genetics, University of British Columbia,
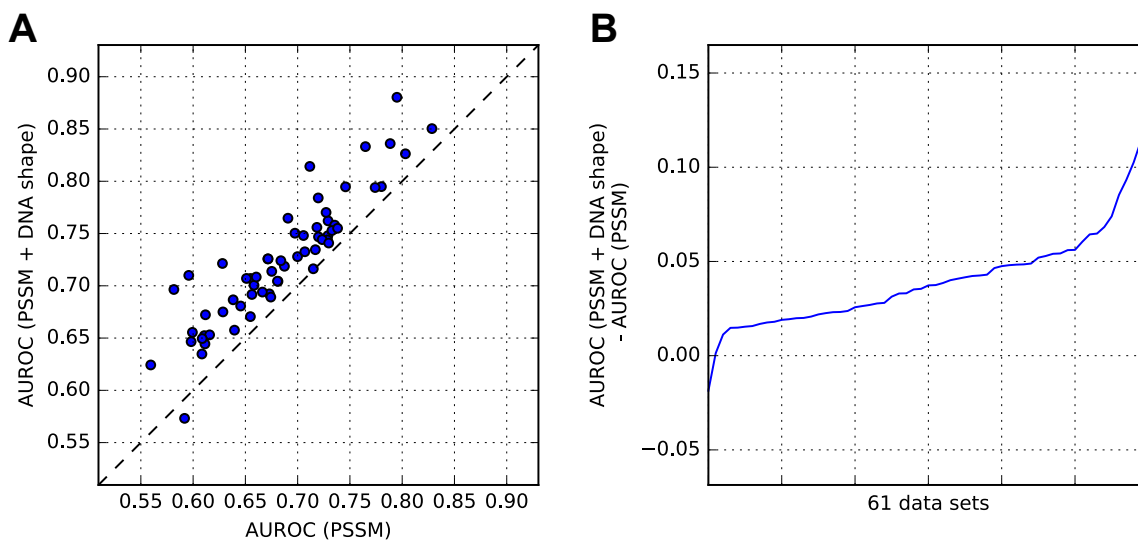980 West 28th Avenue, V5Z 4H4, Vancouver, BC, Canada
[2] Centre for Molecular Medicine Norway (NCMM), Nordic EMBL Partnership,
University of Oslo and Oslo University Hospital, Norway
[3] Molecular and Computational Biology Program, Departments of Biological Sciences,
Chemistry, Physics, and Computer Science,
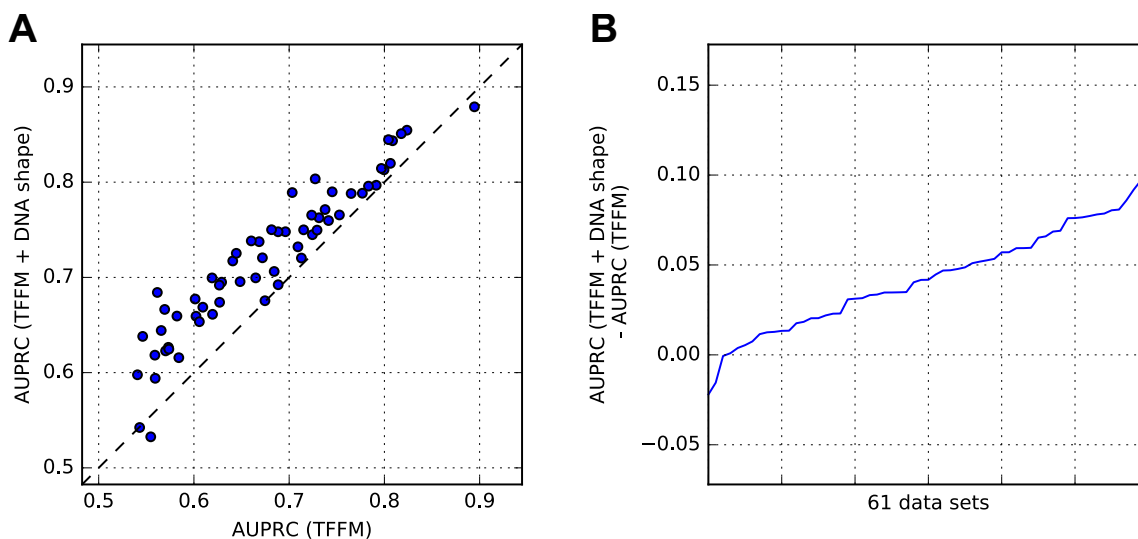University of Southern California, Los Angeles, CA 90089, USA
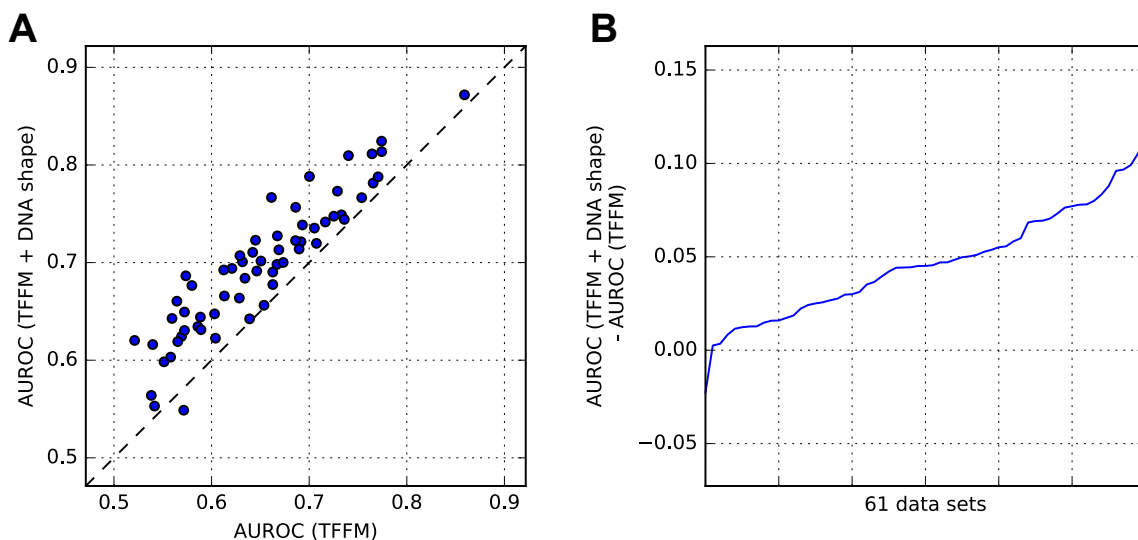* Co-corresponding authors

*Related to Figure 2. Impact of DNA shape on predicting TFBSs when considering recurrent ChIP-seq regions for each TF.* **A.** *Comparison of the AUPRC obtained for the recurrent ENCODE human ChIP-seq peak regions associated with 61 TFs when using either the PSSM scores (x-axis) or the PSSM + DNA shape classifiers (y-axis). The dashed line represents equal AUPRC values obtained with the two methods.* **B.** *Improvement in predictive power obtained when considering DNA shape features (y-axis) as the difference between AUPRC values obtained with the PSSM + DNA shape classifiers and PSSM scores. The higher the difference, the stronger the improvement when DNA shape information is incorporated. Note that the data sets (x-axis) were ranked by increasing difference values.*
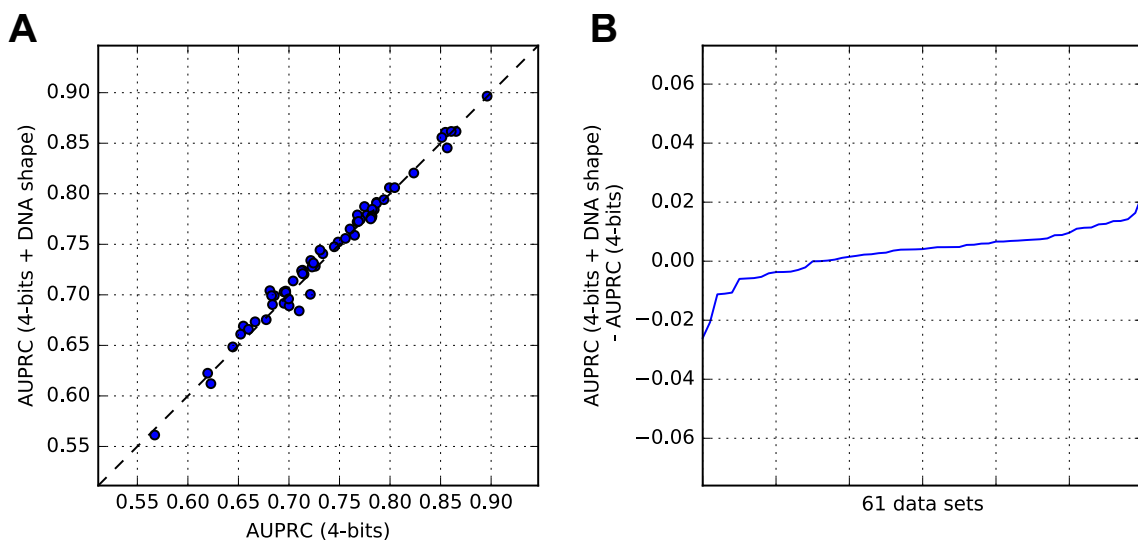
**A**



**B**



*Related to Figure 2. Impact of DNA shape on predicting TFBSs when considering recurrent ChIP-seq regions for each TF. **A.** Comparison of the AUROC obtained for the recurrent ENCODE human ChIP-seq peak regions associated with 61 TFs when using either the PSSM scores (x-axis) or the PSSM + DNA shape classifiers (y-axis). The dashed line represents equal AUROC values obtained with the two methods. **B.** Improvement in predictive power obtained when considering DNA shape features (y-axis) as the difference between AUROC values obtained with the PSSM + DNA shape classifiers and PSSM scores. The higher the difference, the stronger the improvement when DNA shape information is incorporated. Note that the data sets (x-axis) were ranked by increasing difference values.*

**A**



**B**



*Related to Figure 2. Impact of DNA shape on predicting TFBSs when considering recurrent ChIP-seq regions for each TF. **A.** Comparison of the AUPRC obtained for the recurrent ENCODE human ChIP-seq peak regions associated with 61 TFs when using either the TFFM scores (x-axis) or the TFFM + DNA shape classifiers (y-axis). The dashed line represents equal AUPRC values obtained with the two methods. **B.** Improvement in predictive power obtained when considering DNA shape features (y-axis) as the difference between AUPRC values obtained with the TFFM + DNA shape classifiers and TFFM scores. The higher the difference, the stronger the improvement when DNA shape information is incorporated. Note that the data sets (x-axis) were ranked by increasing difference values.*
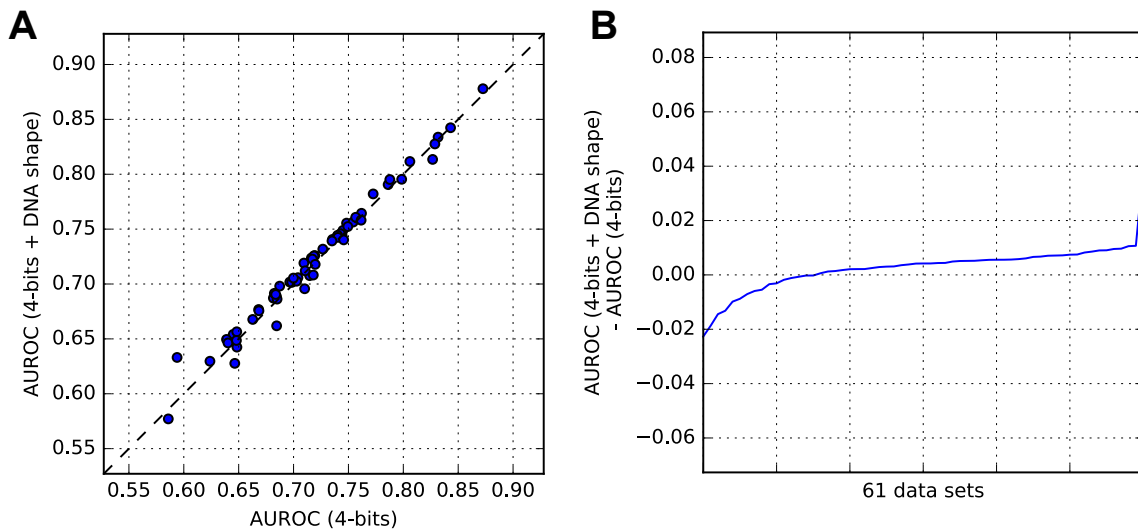
*Related to Figure 2. Impact of DNA shape on predicting TFBSs when considering recurrent ChIP-seq regions for each TF.* **A.** *Comparison of the AUROC obtained for the recurrent ENCODE human ChIP-seq peak regions associated with 61 TFs when using either the TFFM scores (x-axis) or the TFFM + DNA shape classifiers (y-axis). The dashed line represents equal AUROC values obtained with the two methods.* **B.** *improvement in Predictive power obtained when considering DNA shape features (y-axis) as the difference between AUROC values obtained with the TFFM + DNA shape classifiers and TFFM scores. The higher the difference, the stronger the improvement when DNA shape information is incorporated. Note that the data sets (x-axis) were ranked by increasing difference values.*



*Related to Figure 2. Impact of DNA shape on predicting TFBSs when considering recurrent ChIP-seq regions for each TF.* **A.** *Comparison of the AUPRC obtained for the recurrent ENCODE human ChIP-seq peak regions associated with 61 TFs when using either 4-bits (x-axis) or 4-bits + DNA shape (y-axis) classifiers (see (Zhou et al., 2015) for the 4-bits encoding). The dashed line represents equal AUPRC values obtained with the two methods.* **B.** *Improvement in predictive power obtained when considering DNA shape features (y-axis) as the difference between AUPRC values obtained with 4-bits + DNA shape and 4-bits classifiers. The higher the difference, the stronger the improvement when DNA shape information is incorporated. Note that the data sets (x-axis) were ranked by increasing difference values.*

3

*Related to Figure 2. Impact of DNA shape on predicting TFBSs when considering recurrent ChIP-seq regions for each TF. **A.** Comparison of the AUROC obtained for the recurrent ENCODE human ChIP-seq peak regions associated with 61 TFs when using either 4-bits (x-axis) or 4-bits + DNA shape (y-axis) classifiers (see (Zhou et al., 2015) for the 4-bits encoding). The dashed line represents equal AUROC values obtained with the two methods. **B.** Improvement in predictive power obtained when considering DNA shape features (y-axis) as the difference between AUROC values obtained with 4-bits + DNA shape and 4-bits classifiers. The higher the difference, the stronger the improvement when DNA shape information is incorporated. Note that the data sets (x-axis) were ranked by increasing difference values.*

4