

---

## Data S7

# DNA shape features improve transcription factor binding site prediction *in vivo*

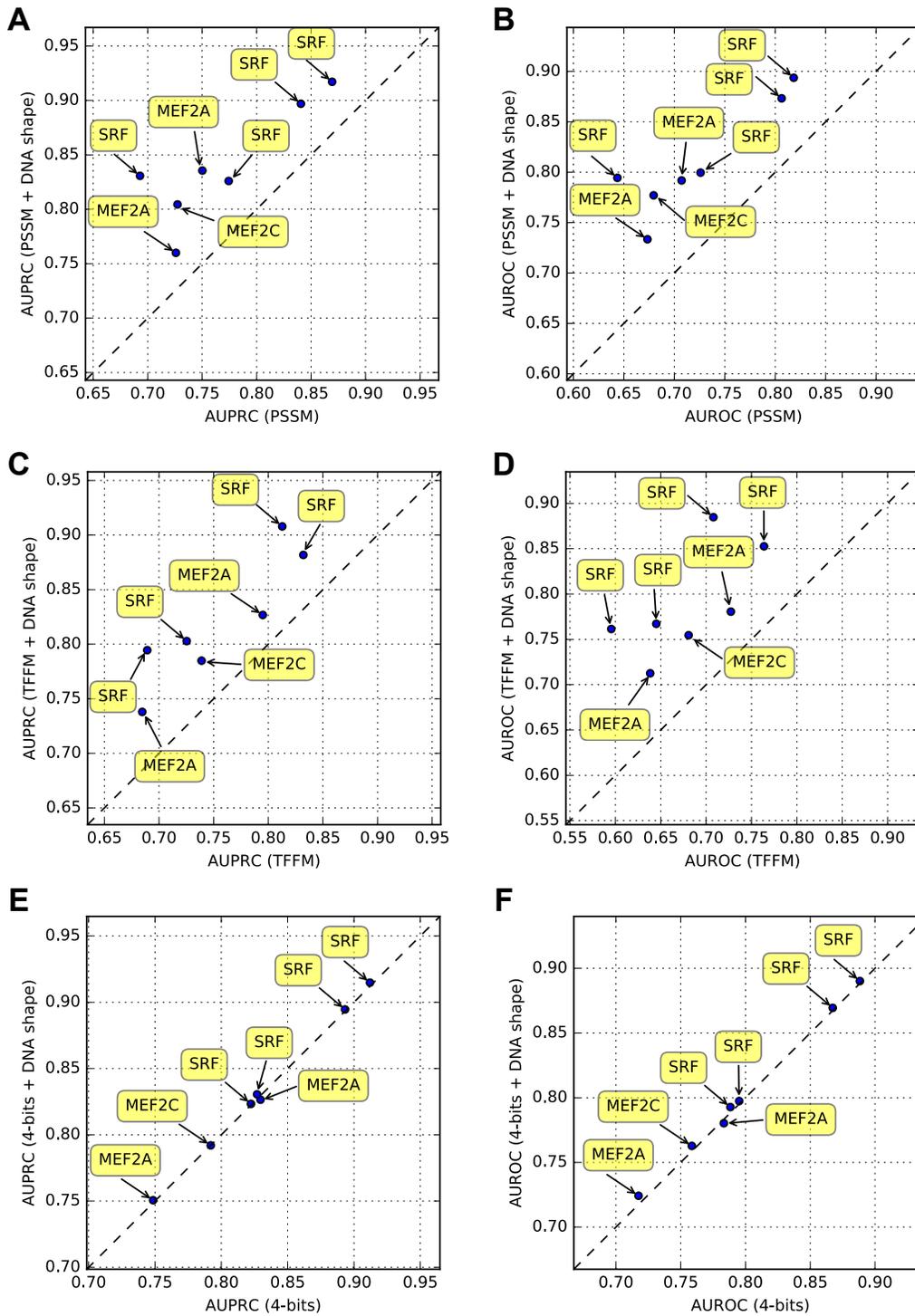
ANTHONY MATHÉLIER<sup>1,2</sup>, BEIBEI XIN<sup>3</sup>, TSU-PEI CHIU<sup>3</sup>,  
LIN YANG<sup>3</sup>, REMO ROHS<sup>3,\*</sup>, AND WYETH W. WASSERMAN<sup>1,\*</sup>

<sup>1</sup> 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

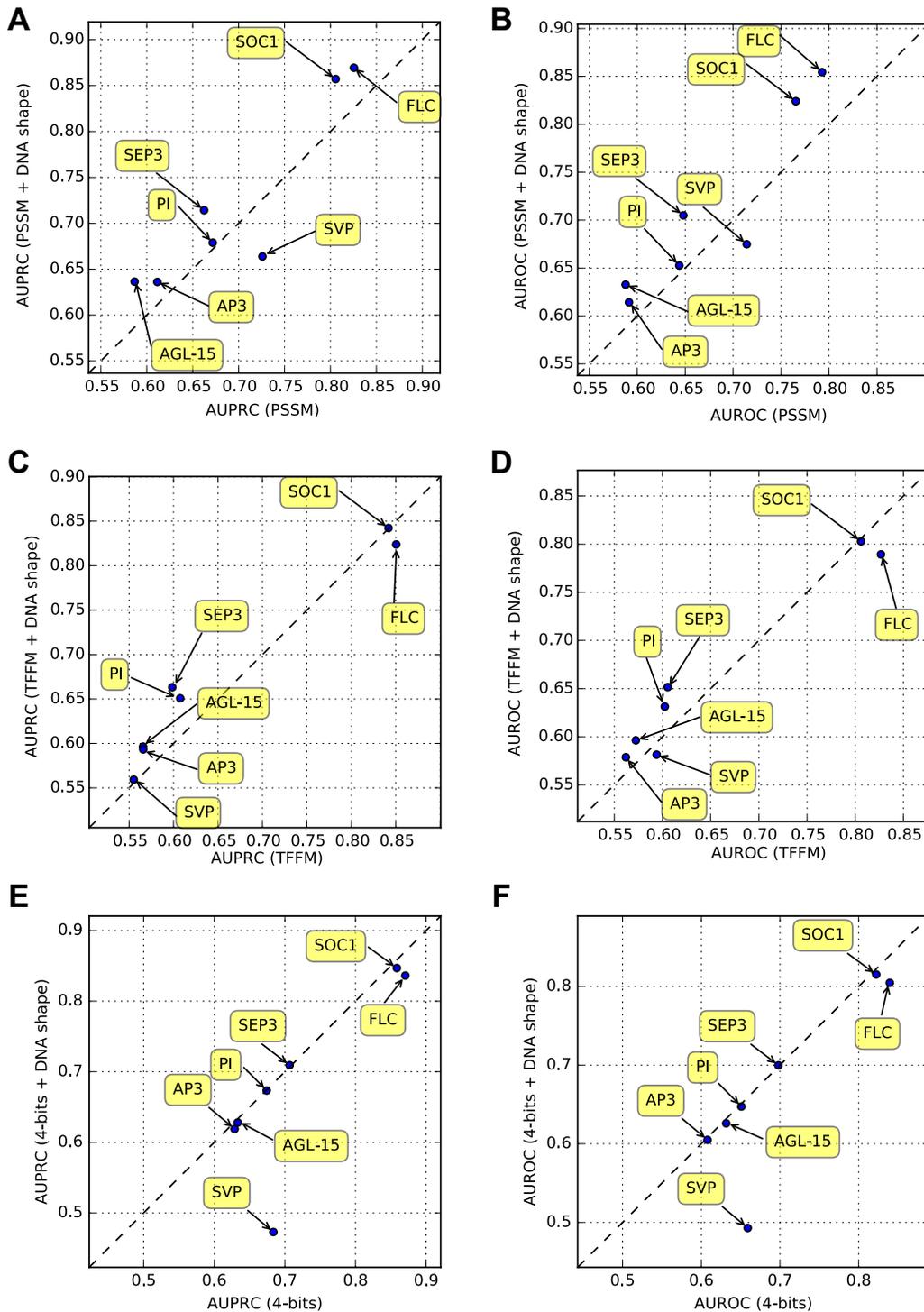
<sup>2</sup> Centre for Molecular Medicine Norway (NCMM), Nordic EMBL partnership,  
University of Oslo and Oslo University Hospital, Norway

<sup>3</sup> Molecular and Computational Biology Program, Departments of Biological Sciences,  
Chemistry, Physics, and Computer Science,  
University of Southern California, Los Angeles, CA 90089, USA

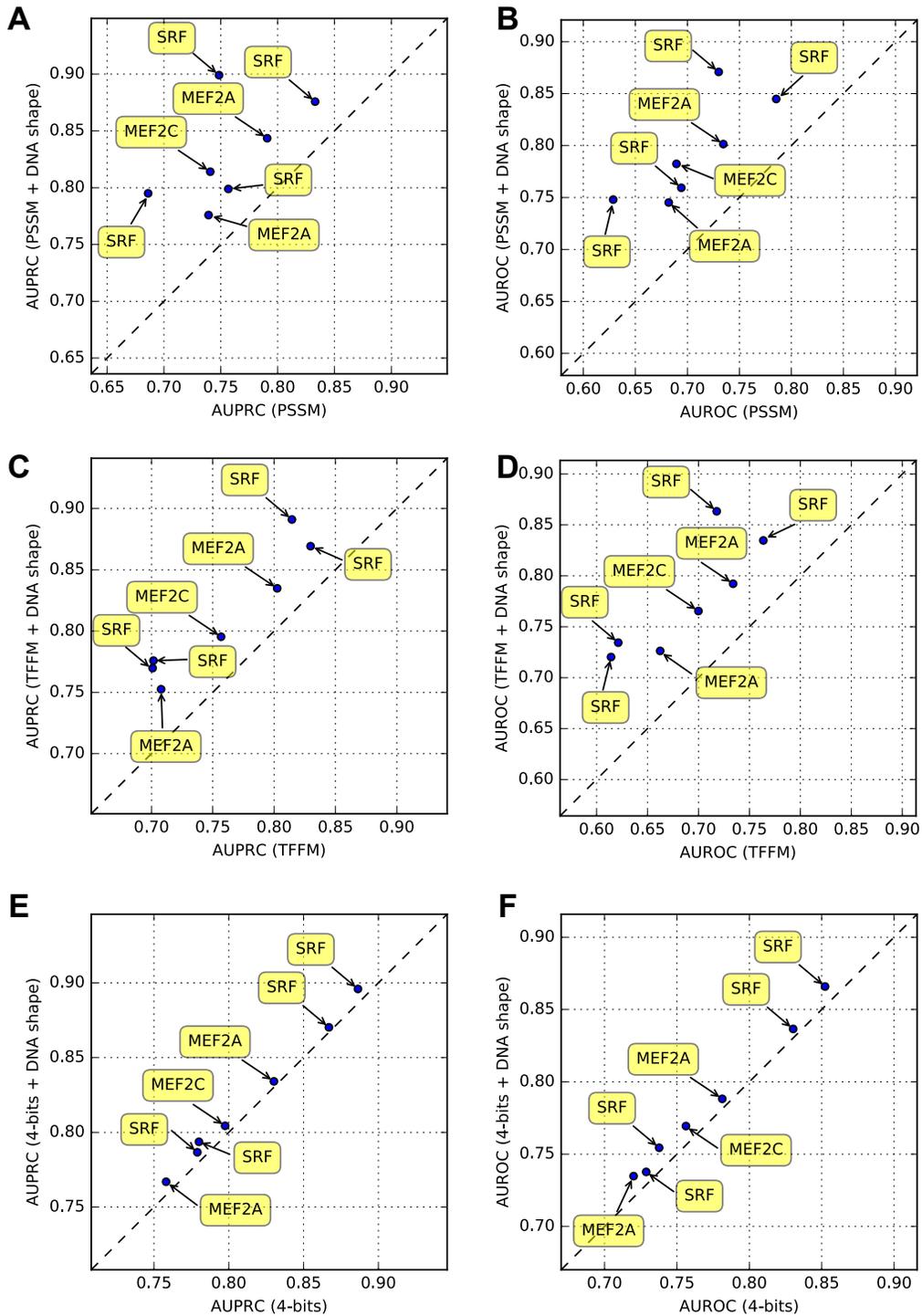
\* Co-corresponding authors



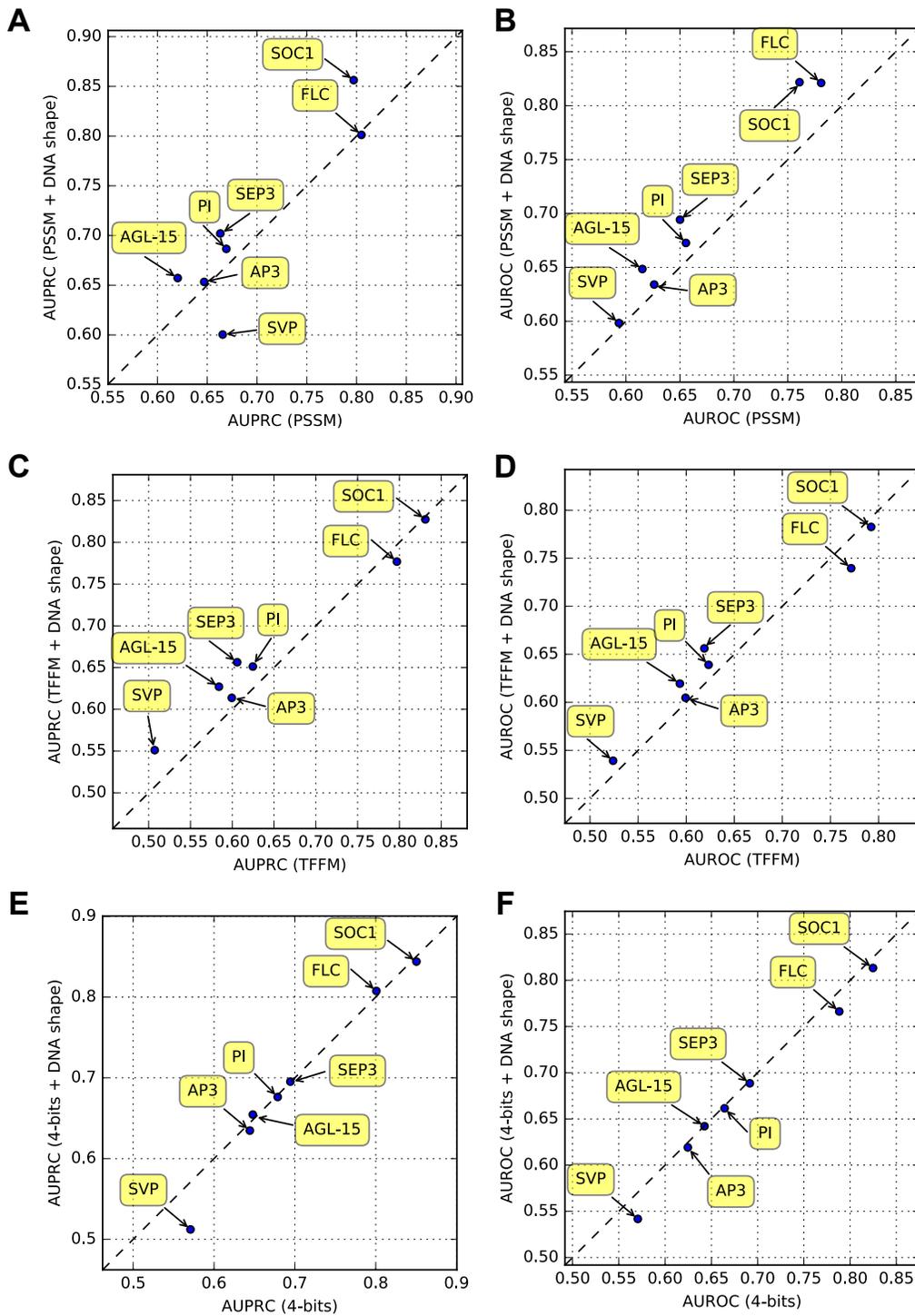
Related to Figure 2. Comparison of the AUPRC (A, C, and E) or AUROC (B, D, and F) obtained for the 7 human ENCODE MADS TF ChIP-seq data sets when using the PSSM scores (x-axis; A-B), the TFFM scores (x-axis; C-D), or the 4-bits classifier (x-axis; E-F; see (Zhou et al., 2015) for the 4-bits encoding) versus the PSSM + DNA shape (y-axis; A-B), TFFM + DNA shape (y-axis; C-D), or 4-bits + DNA shape (y-axis; E-F) classifiers.



Related to Figure 2. Comparison of the AUPRC (A, C, and E) or AUROC (B, D, and F) obtained for the 7 plant MADS TF ChIP-seq data sets when using the PSSM scores (x-axis; A-B), the TFFM scores (x-axis; C-D), or the 4-bits classifier (x-axis; E-F; see (Zhou et al., 2015) for the 4-bits encoding) versus the PSSM + DNA shape (y-axis; A-B), TFFM + DNA shape (y-axis; C-D), or 4-bits + DNA shape (y-axis; E-F) classifiers.



Related to Figure 2. Impact of DNA shape on predicting human MADS-box TFBSs with background sequences matching the dinucleotide composition of ChIP-seq regions. Comparison of the AUPRC (A, C, and E) or AUROC (B, D, and F) obtained for the 7 human ENCODE MADS TF ChIP-seq data sets when using the PSSM scores (x-axis; A-B), the TFFM scores (x-axis; C-D), or the 4-bits classifier (x-axis; E-F; see (Zhou et al., 2015) for the 4-bits encoding) versus the PSSM + DNA shape (y-axis; A-B), TFFM + DNA shape (y-axis; C-D), or 4-bits + DNA shape (y-axis; E-F) classifiers.



Related to Figure 2. Impact of DNA shape on predicting plant MADS-box TFBSs with background sequences matching the dinucleotide composition of ChIP-seq regions. Comparison of the AUPRC (A, C, and E) or AUROC (B, D, and F) obtained for the 7 plant MADS TF ChIP-seq data sets when using the PSSM scores (x-axis; A-B), the TFFM scores (x-axis; C-D), or the 4-bits classifier (x-axis; E-F; see (Zhou et al., 2015) for the 4-bits encoding) versus the PSSM + DNA shape (y-axis; A-B), TFFM + DNA shape (y-axis; C-D), or 4-bits + DNA shape (y-axis; E-F) classifiers.