**Supplemental information**

# DeepFace: Deep-learning-based framework

# to contextualize orofacial-cleft-related variants

# during human embryonic craniofacial development

Yulin Dai, Toshiyuki Itai, Guangsheng Pei, Fangfang Yan, Yan Chu, Xiaoqian Jiang, Seth M. Weinberg, Nandita Mukhopadhyay, Mary L. Marazita, Lukas M. Simon, Peilin Jia, and Zhongming Zhao
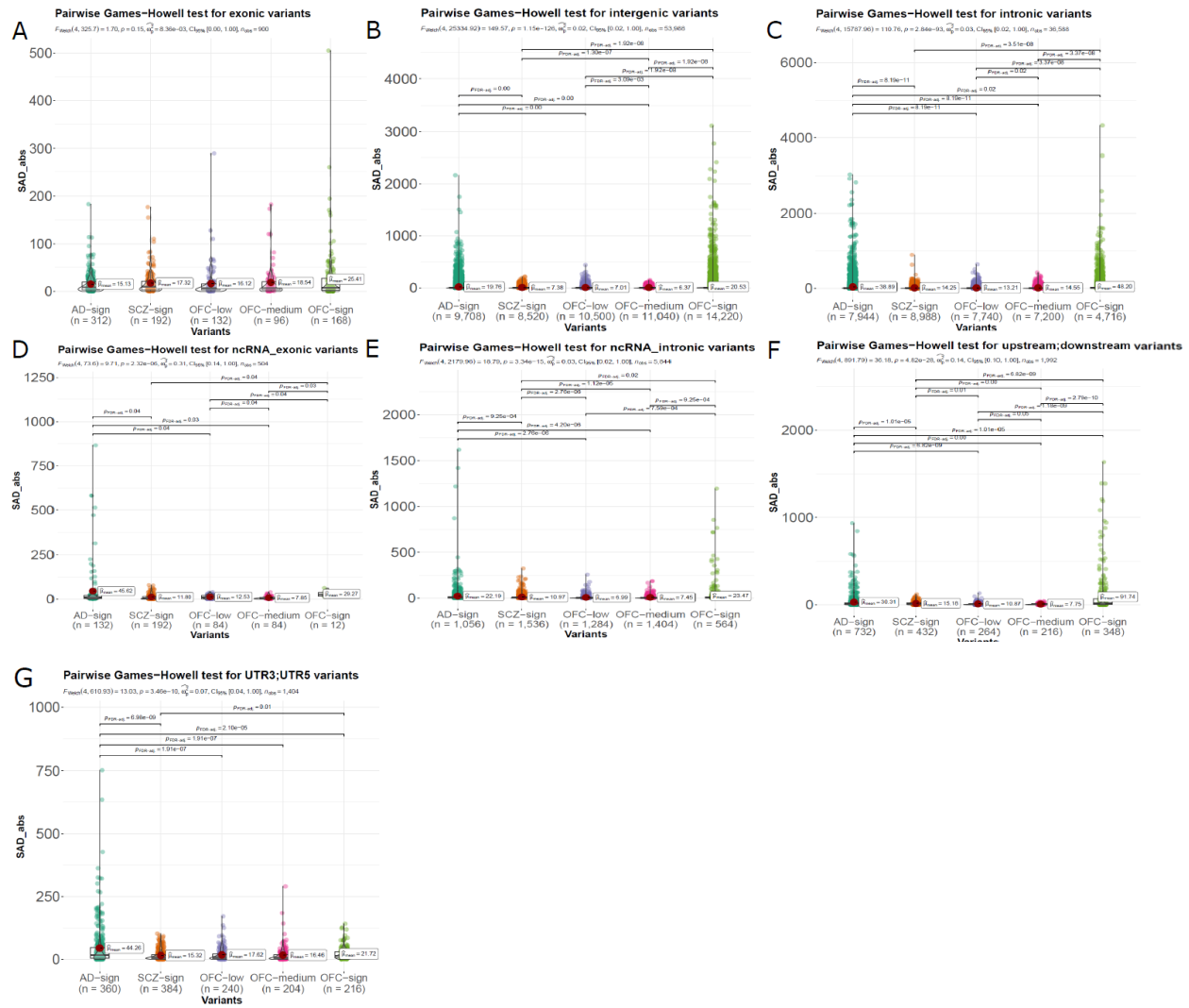
**Figure S1. Pairwise Games-Howell test for seven variant categories across five different SNP sets**

Panels (A) to (G) depict the results of Welch's F-test from ANOVA, conducted to determine whether there is a mean difference across comparison groups. Subsequently, the non-parametric Games-Howell test was used to evaluate whether pairwise mean rank differences exist in the absolute SAD scores between the five variant categories: AD-sign, SCZ-sign, OFC-low, OFC-medium, and OFC-sign. In each sub-panel, only the comparison groups with significant $P_{FDR}$-adjusted p-values are highlighted.
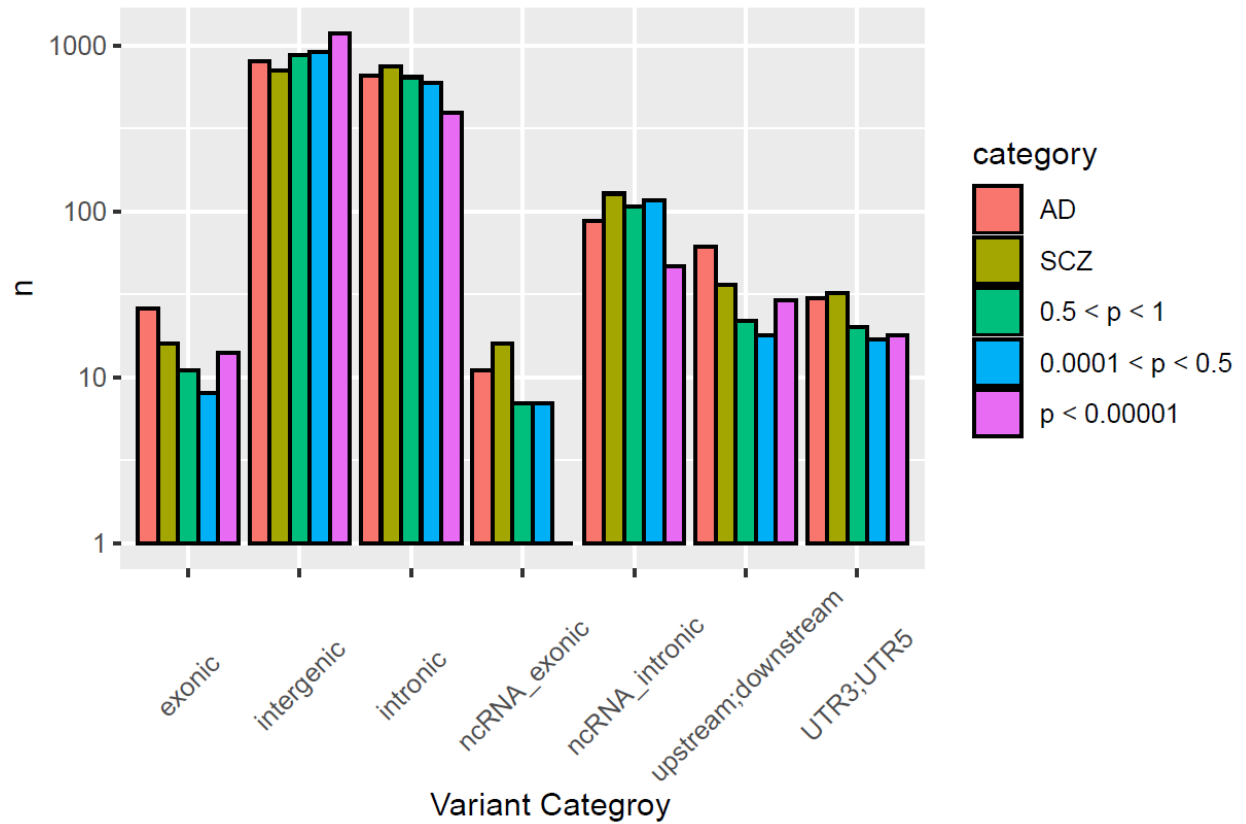
**Figure S2. Number of SNP by variant category for five different SNP sets**

Number of variants from five different categories AD-sign, SCZ-sign, OFC-low, OFC-medium, and OFC-sign stratified by variant category.
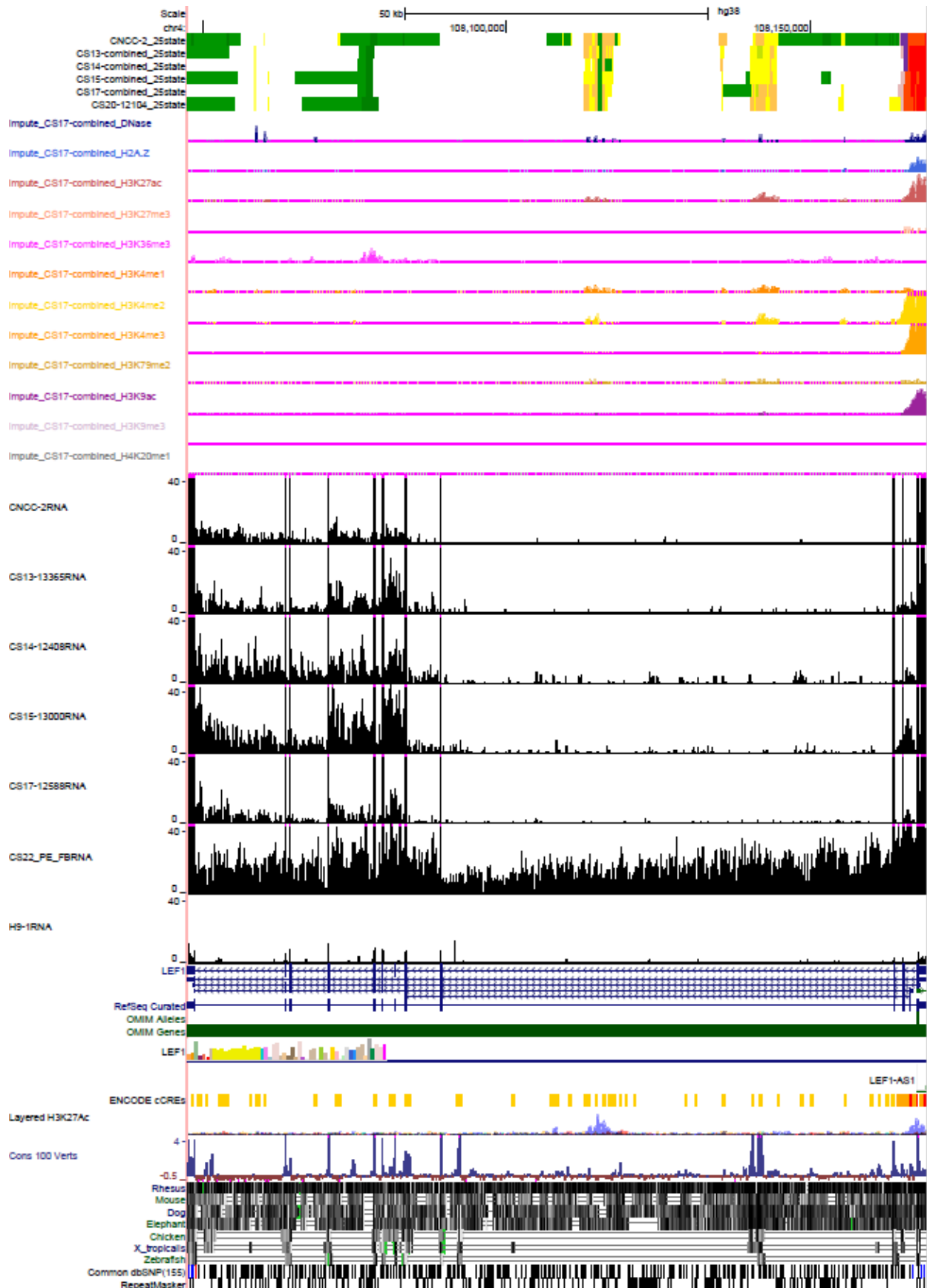
**Figure S3. LEF1 in UCSC from Cotney lab craniofacial Genome Browser for transcriptome and epigenomic features for craniofacial tissue from Carnegie stages CS 13 to CS 22.**
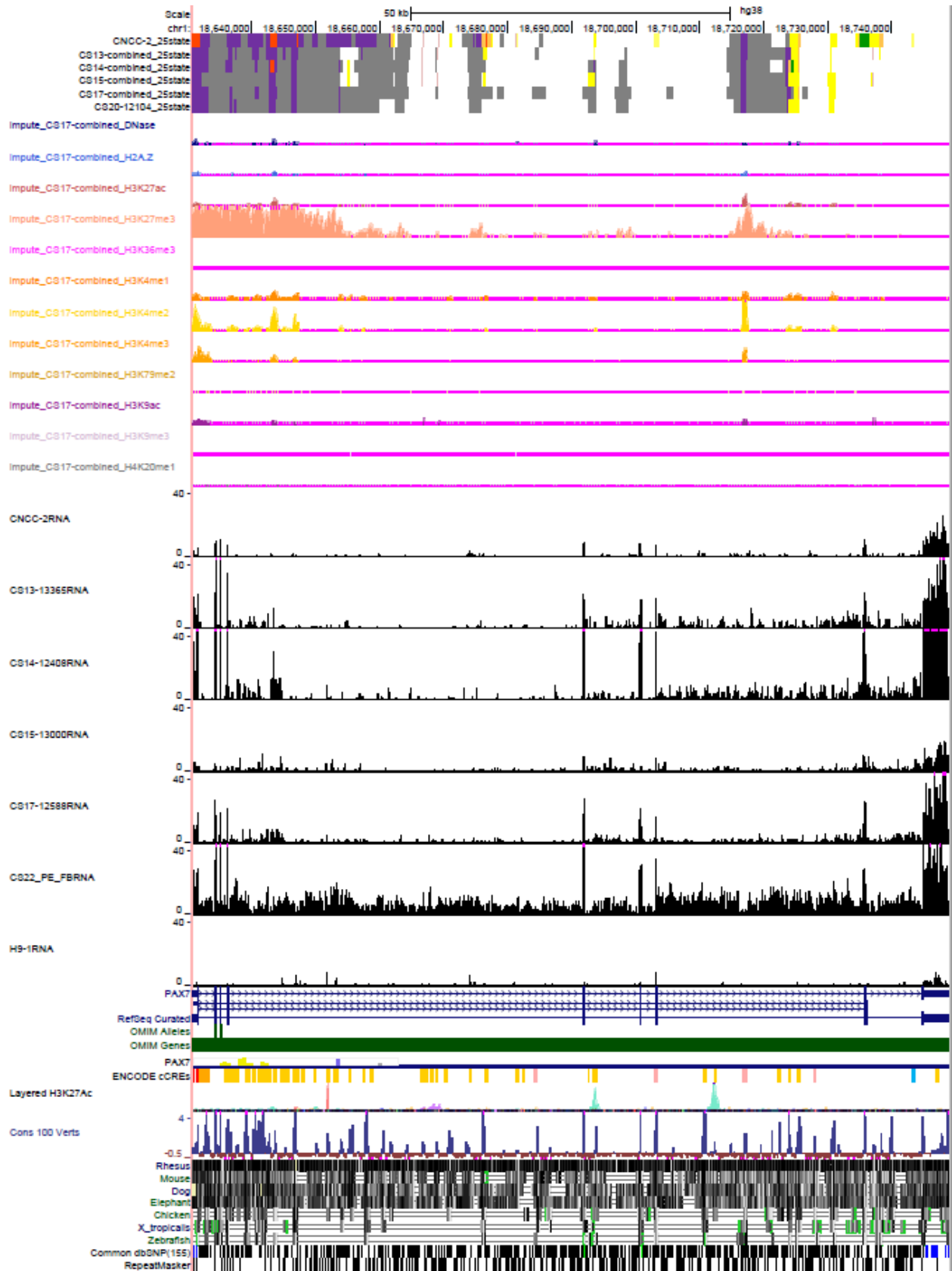
**Figure S4. PAX7 in UCSC from Cotney lab craniofacial Genome Browser for transcriptome and epigenomic features for craniofacial tissue from Carnegie stages CS 13 to CS 22.**
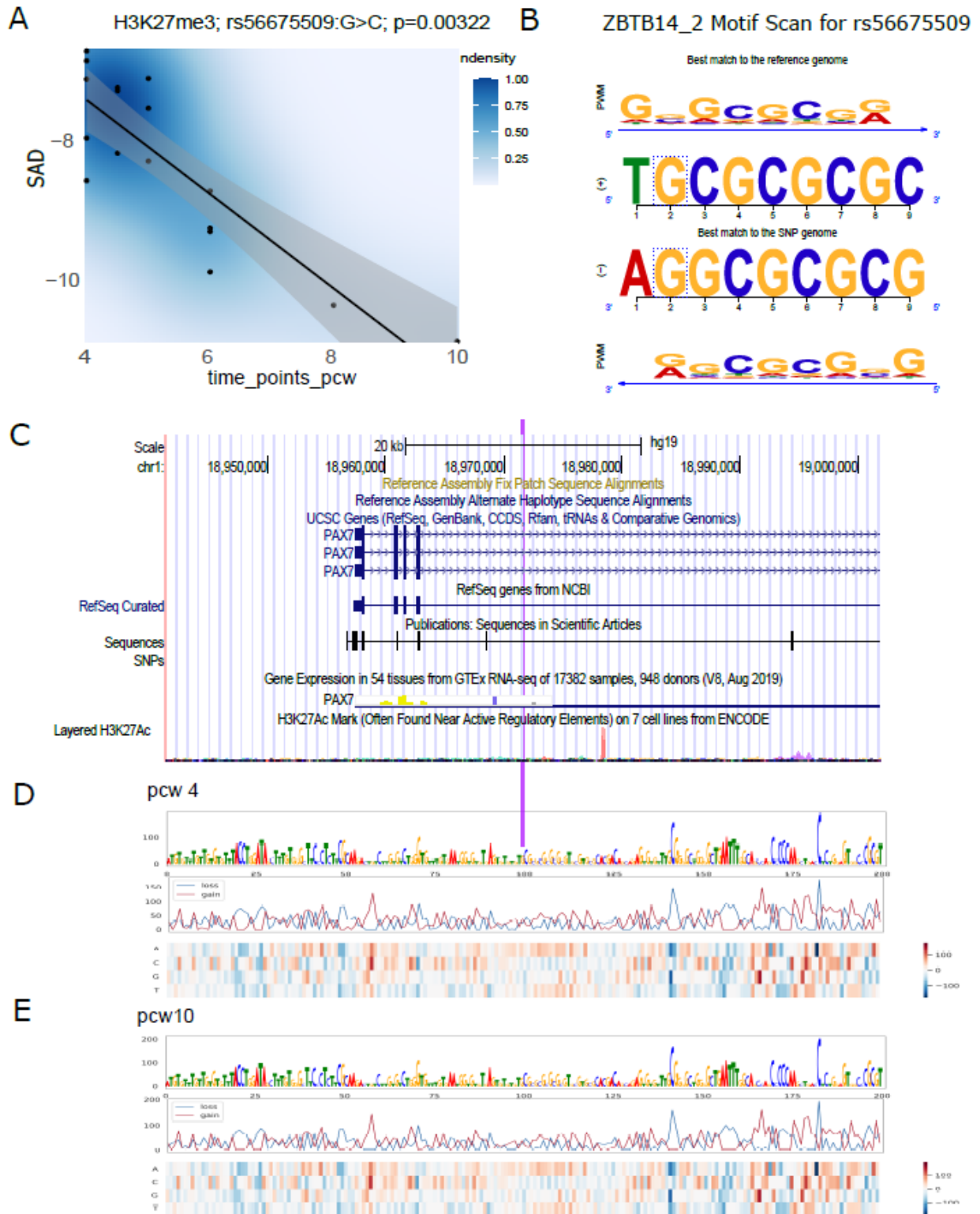
**Figure S5. Motif analysis of DeepFace variant rs56675509.**

(A) SNP activity difference (SAD) for rs56675509 (G>C) associated with post-conception weeks (pcw) increase for H3K27me3. The blue area of intensity is positively correlated with the SAD score point density. The black line and grey confidence interval (95%) model the linear relationship between SAD scores and pcw. P-value indicates the possibility of null hypothesis that the coefficient is equal to zero is true. (B) Sequence logo stacks from top to bottom: sequence logo of reference allele matching position weight matrix; Reference subsequences, alternative allele subsequences, and sequence logo of SNP allele matching position weight matrix. The best match reference sequence and alternative allele sequence for the motif ZBTB14 were visualized. (C) UCSC genome browser for the rs56675509 and its surrounding gene. The purple vertical line indicates the exact genomic region of 200 bps for (D) and (E). (D) & (E) show the dynamic gain and loss of SAD score for all possible substitutions in each of the 200-bp genomic positions around the rs56675509 in pcw 4 and pcw 10, respectively. The alteration between D & E is relatively small in figure. These SAD score dynamics were visualized in three ways: 1) sequence logo weight by the loss of SAD across 200-bp sequence; 2) The blue and red lines indicate the minimum (loss) and maximum (gain) change among the possible substitutions from reference allele; 3) The quantities in the heatmap display the change in SAD after substituting nucleotide from reference allele.
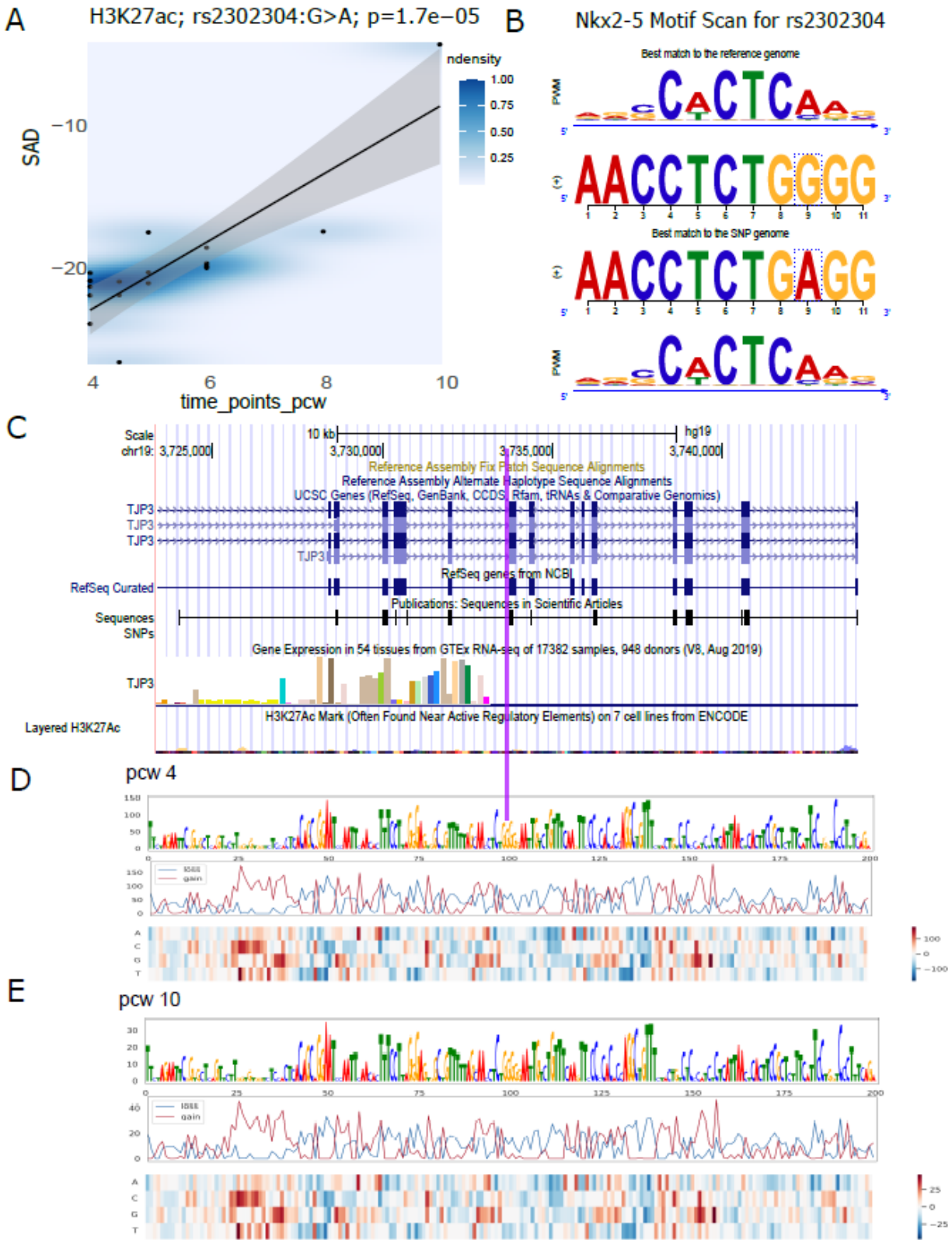
**Figure S6. Motif analysis of DeepFace variant rs2302304.**

(A) SNP activity difference (SAD) for rs2302304 (G>A) associated with post-conception weeks (pcw) increase for H3K27ac. The blue area of intensity is positively correlated with the SAD score point density. The black line and grey confidence interval (95%) model the linear relationship between SAD scores and pcw. P-value indicates the possibility of null hypothesis that the coefficient is equal to zero is true. (B) Sequence logo stacks from top to bottom: sequence logo of reference allele matching position weight matrix; Reference subsequences, alternative allele subsequences, and sequence logo of SNP allele matching position weight matrix. The best match reference sequence and alternative allele sequence for the motif Nkx2-5 were visualized (C) UCSC genome browser for the rs2302304 and its surrounding gene. The purple vertical line indicates the exact genomic region of 200 bps for (D) and (E). (D) & (E) show the dynamic gain and loss of SAD score for all possible substitutions in each of the 200-bp genomic positions around the rs2302304 in pcw 4 and pcw 10, respectively. The alteration between D & E is relatively small in figure.These SAD score dynamics were visualized in three ways: 1) sequence logo weight by the loss of SAD across 200-bp sequence; 2) The blue and red lines indicate the minimum (loss) and maximum (gain) change among the possible substitutions from reference allele; 3) The quantities in the heatmap display the change in SAD after substituting nucleotide from reference allele.
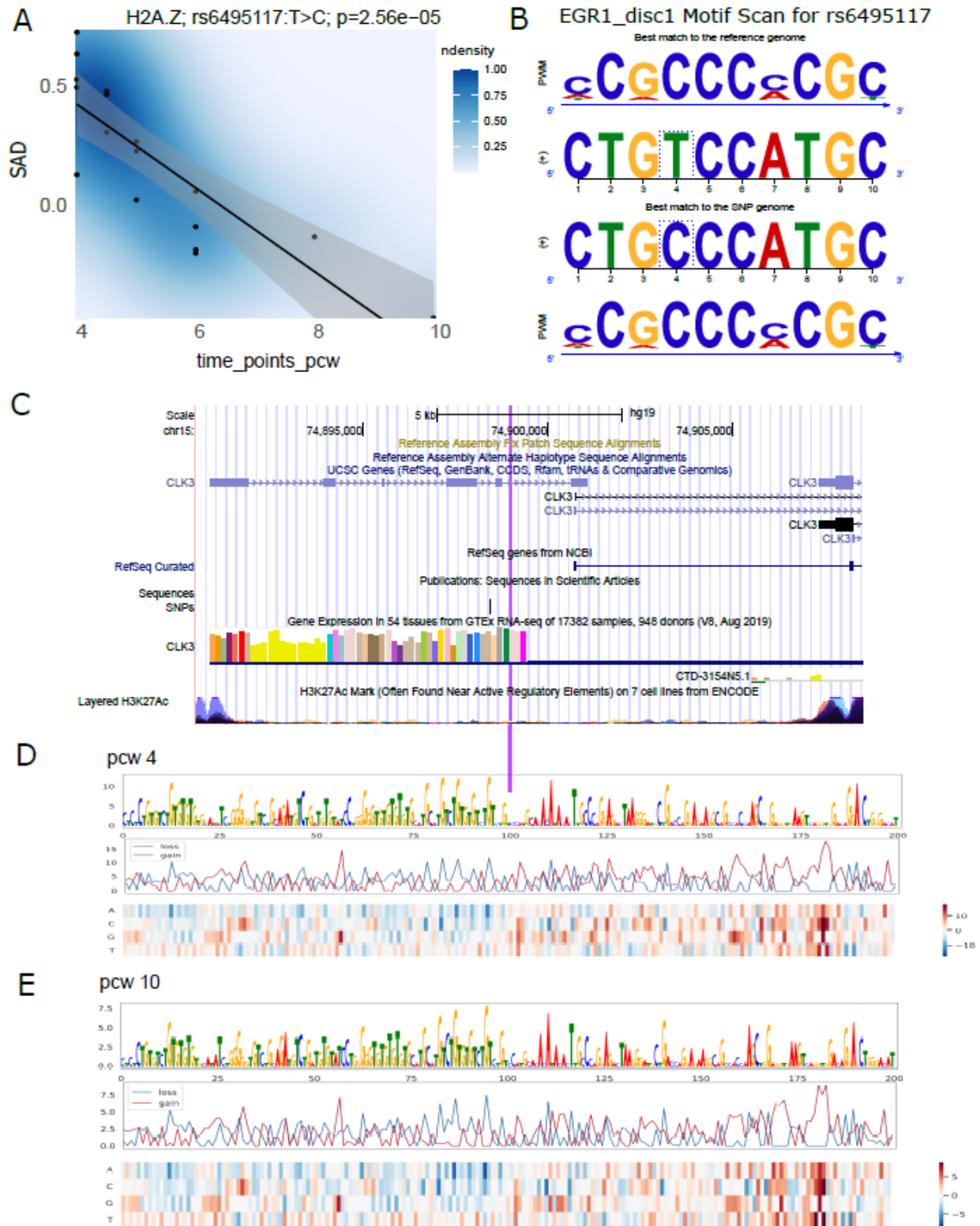
**Figure S7. Motif analysis of DeepFace variant rs6495117.**

(A) SNP activity difference (SAD) for rs6495117 (T>C) associated with post-conception weeks (pcw) increase for H2A.Z. The blue area of intensity is positively correlated with the SAD score point density. The black line and grey confidence interval (95%) model the linear relationship between SAD scores and pcw. P-value indicates the possibility of null hypothesis that the coefficient is equal to zero is true. (B) Sequence logo stacks from top to bottom: sequence logo of reference allele matching position weight matrix; Reference subsequences, alternative allele subsequences, and sequence logo of SNP allele matching position weight matrix. The best match reference sequence and alternative allele sequence for the motif EGR1 were visualized (C) UCSC genome browser for the rs6495117 and its surrounding gene. The purple vertical line indicates the exact genomic region of 200 bps for (D) and (E). (D) & (E) show the dynamic gain and loss of SAD score for all possible substitutions in each of the 200-bp genomic positions around the rs6495117 in pcw 4 and pcw 10, respectively. The alteration between D & E is relatively small in figure. These SAD score dynamics were visualized in three ways: 1) sequence logo weight by the loss of SAD across 200-bp sequence; 2) The blue and red lines indicate the minimum (loss) and maximum (gain) change among the possible substitutions from reference allele; 3) The quantities in the heatmap display the change in SAD after substituting nucleotide from reference allele.
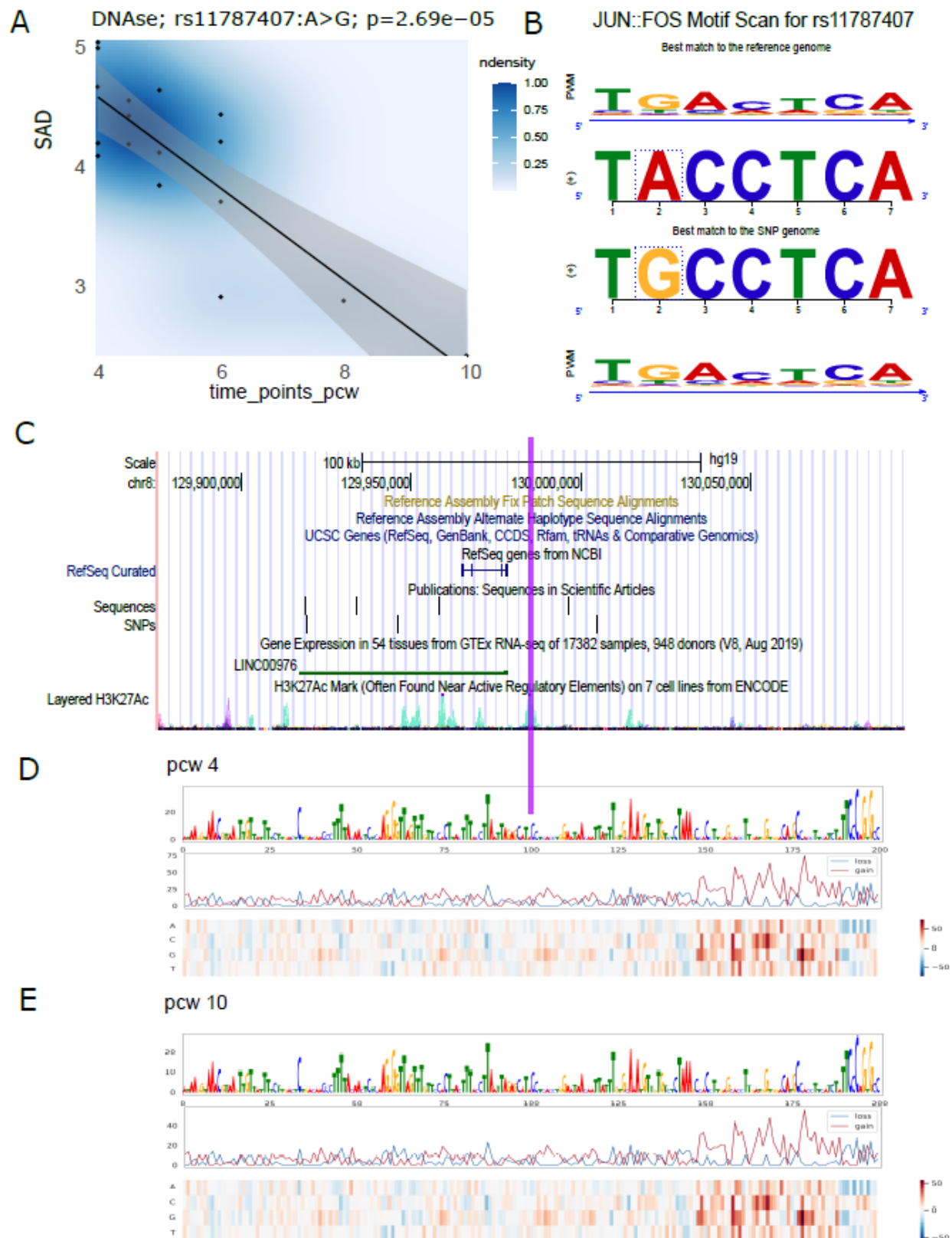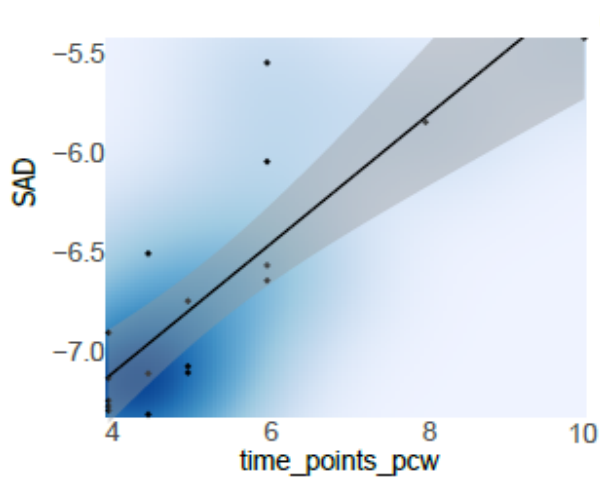
**Figure S8. Motif analysis of DeepFace variant rs11787407.**

(A) SNP activity difference (SAD) for rs11787407 (A>G) associated with post-conception weeks (pcw) increase for DNase. The blue area of intensity is positively correlated with the SAD score point density. The black line and grey confidence interval (95%) model the linear relationship between SAD scores and pcw. P-value indicates the possibility of null hypothesis that the coefficient is equal to zero is true. (B) Sequence logo stacks from top to bottom: sequence logo of reference allele matching position weight matrix; Reference subsequences, alternative allele subsequences, and sequence logo of SNP allele matching position weight matrix. The best match reference sequence and alternative allele sequence for the motif JUN/FOS were visualized (C) UCSC genome browser for the rs11787407 and its surrounding gene. The purple vertical line indicates the exact genomic region of 200 bps for (D) and (E). (D) & (E) show the dynamic gain and loss of SAD score for all possible substitutions in each of the 200-bp genomic positions around the rs11787407 in pcw 4 and pcw 10, respectively. These SAD score dynamics were visualized in three ways: 1) sequence logo weight by the loss of SAD across 200-bp sequence; 2) The blue and red lines indicate the minimum (loss) and maximum (gain) change among the possible substitutions from reference allele; 3) The quantities in the heatmap display the change in SAD after substituting nucleotide from reference allele.
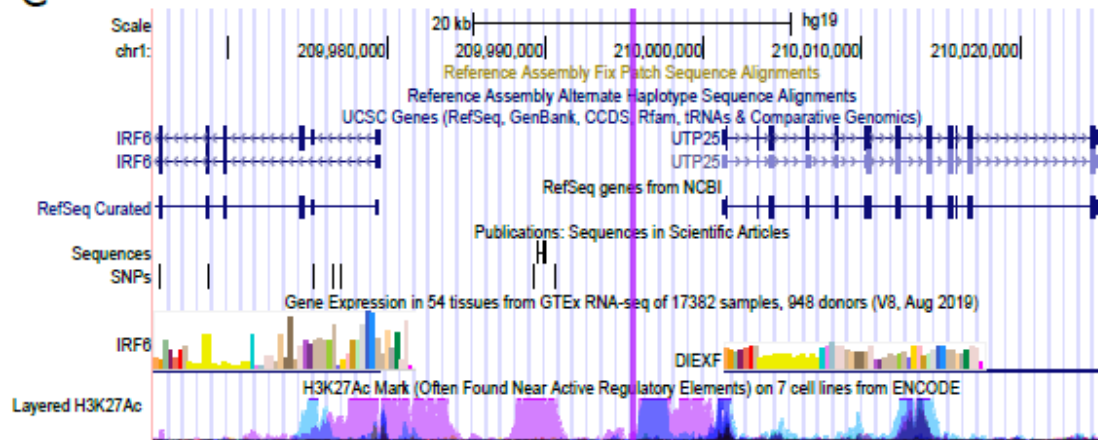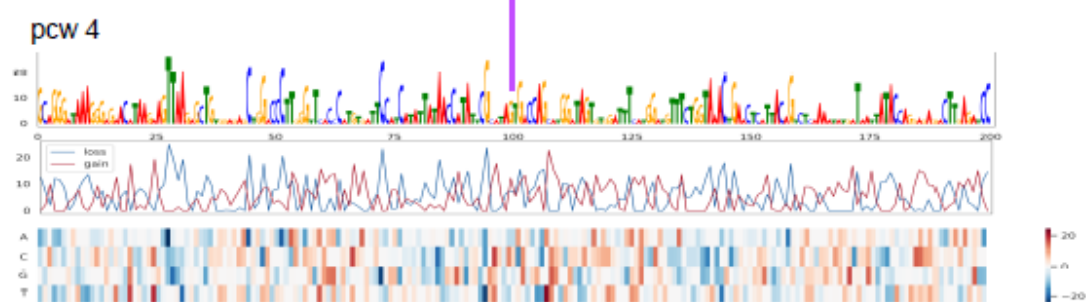
**A** DNAse; rs12075674:G>A; p=1.94e−05

**B** AFP_1 Motif Scan for rs12075674
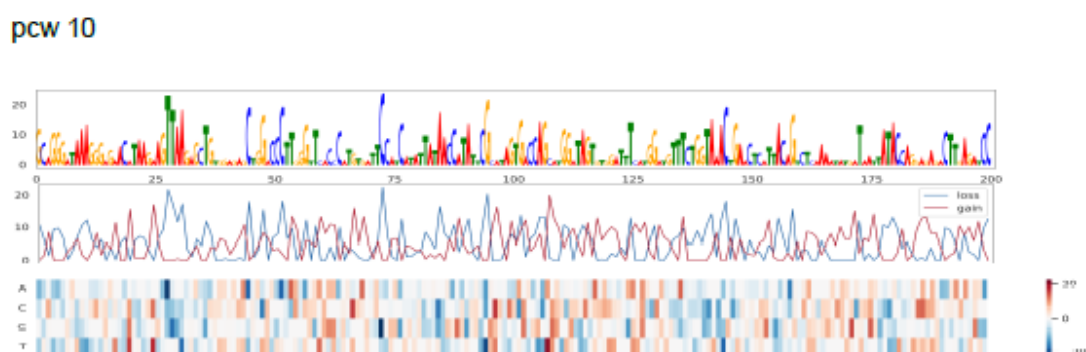
**C**

**D** pcw 4

**E** pcw 10

**Figure S9. Motif analysis of DeepFace variant rs12075674.**

(A) SNP activity difference (SAD) for rs12075674 (G>A) associated with post-conception weeks (pcw) increase for DNase. The blue area of intensity is positively correlated with the SAD score point density. The black line and grey confidence interval (95%) model the linear relationship between SAD scores and pcw. P-value indicates the possibility of null hypothesis that the coefficient is equal to zero is true. (B) Sequence logo stacks from top to bottom: sequence logo of reference allele matching position weight matrix; Reference subsequences, alternative allele subsequences, and sequence logo of SNP allele matching position weight matrix. The best match reference sequence and alternative allele sequence for the motif AFP were visualized (C) UCSC genome browser for the rs12075674 and its surrounding gene. The purple vertical line indicates the exact genomic region of 200 bps for (D) and (E). (D) & (E) show the dynamic gain and loss of SAD score for all possible substitutions in each of the 200-bp genomic positions around the rs12075674 in pcw 4 and pcw 10, respectively. These SAD score dynamics were visualized in three ways: 1) sequence logo weight by the loss of SAD across 200-bp sequence; 2) The blue and red lines indicate the minimum (loss) and maximum (gain) change among the possible substitutions from reference allele; 3) The quantities in the heatmap display the change in SAD after substituting nucleotide from reference allele.