



# CALIFORNIA STATE SCIENCE FAIR 2015 PROJECT SUMMARY

<b>Name(s)</b> <b>Simon L. Kuang</b>	<b>Project Number</b> <b>S1414</b>
<b>Project Title</b> <b>Predicting Chromosome Conformation from Epigenetic Features</b>	
<b>Abstract</b> <b>Objectives/Goals</b> The objective is to predict the results of 40kb <i>HindIII</i> Hi-C experimentation on the mouse ESC genome from 21 documented mouse epigenetic features. <b>Methods/Materials</b> Around 200 million 22-dimensional datapoints were drawn from the NIH histone database and from the Hi-C paper. Characteristics of the data required development of a novel machine learning meta-algorithm using inverse probability Metropolis-Hastings sampling, bagging with 100 independent samples, single-layer perceptrons optimized using Levenberg-Marquardt training, and intelligent weighting based on in-sample prediction correlation. Chromosome 1 was used for training, chromosome 2 for validation, and chromosomes 3 to X for testing. <b>Results</b> The model's predictions scored a Pearson correlation coefficient between 0.7 and 0.8 on all chromosomes. In-sample RMSE and out-of-sample RMSE decreased on all predictors over the 5 training epochs, and in-sample correlation and out-of-sample correlation increased on all predictors over the 5 training epochs. <b>Conclusions/Discussion</b> My method performs extremely well on this unusual dataset. In the next decade, spatial distance will be more insightful, cheaper and faster to compute, and easier to analyze in clinical situations than the DNA sequence itself.	
<b>Summary Statement</b> A novel machine learning meta-algorithm is necessary and effective for studying the 3-dimensional conformation of the mammalian genome.	
<b>Help Received</b> Used computing facilities at PSU Center for Comparative Genomics and Bioinformatics	