

A All Consistency Plots

GPSR (high variability setting) uses a small population (20 individuals) evolved over many generations (200) with low operator probabilities, $p_{\text{crossover}} = 0.09$, $p_{\text{mutation}} = 0.001$, while GPSR (low variability setting) employs a large population (200 individuals) evolved over few generations (20) with high operator probabilities $p_{\text{crossover}} = 0.9$, $p_{\text{mutation}} = 0.01$.

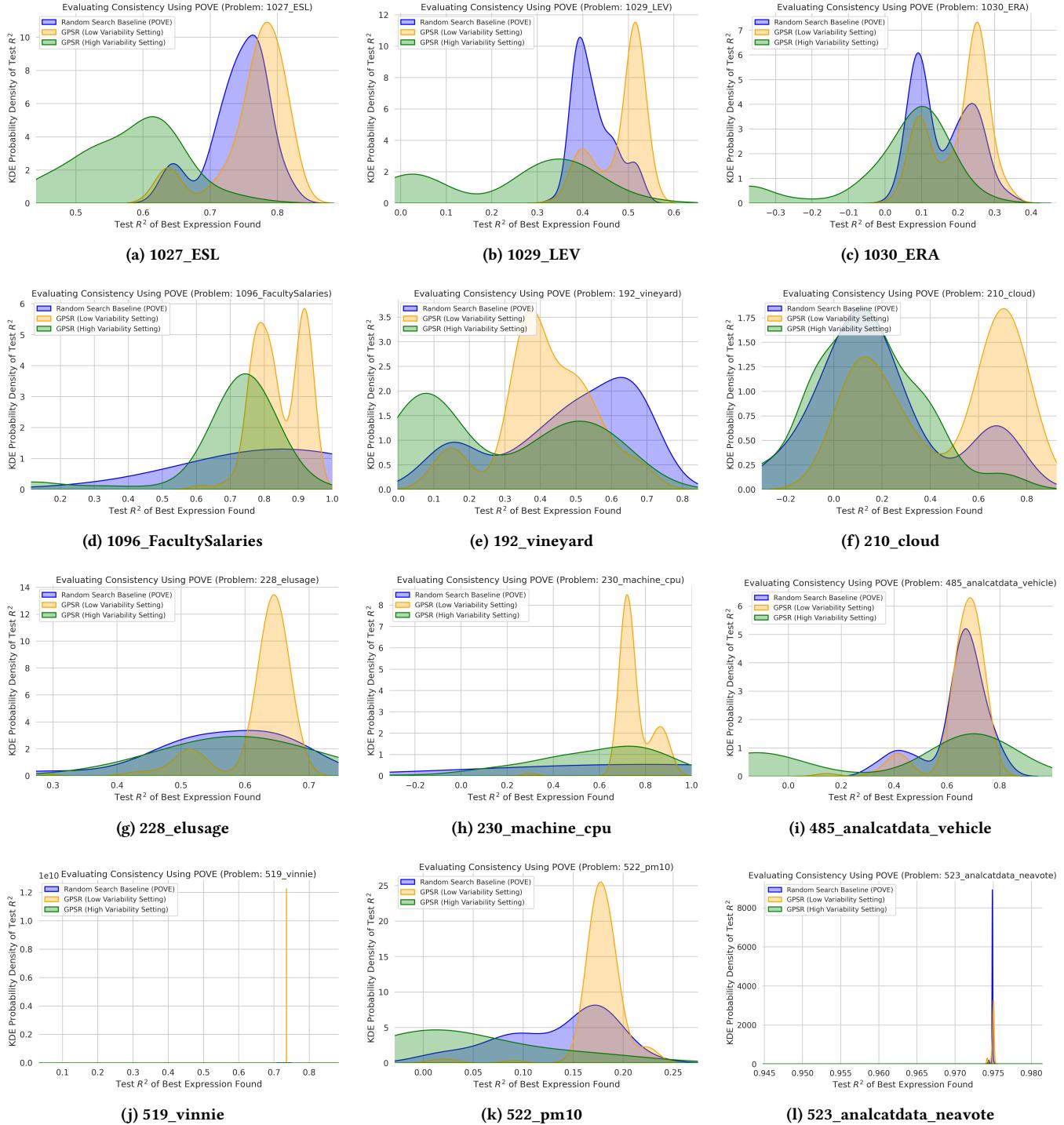


Figure 1: 1st to 12th problem

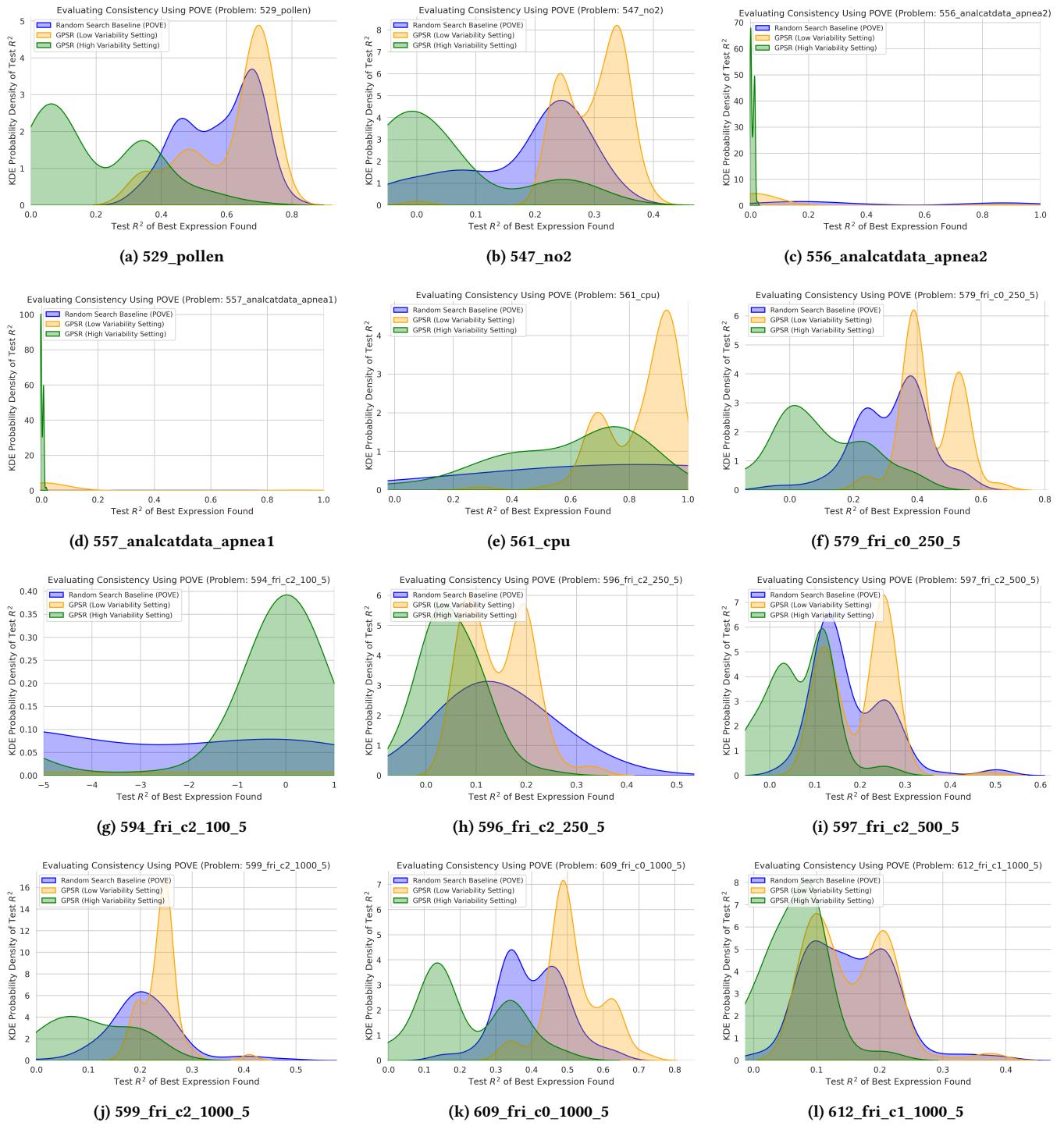


Figure 2: 13th to 24th problem

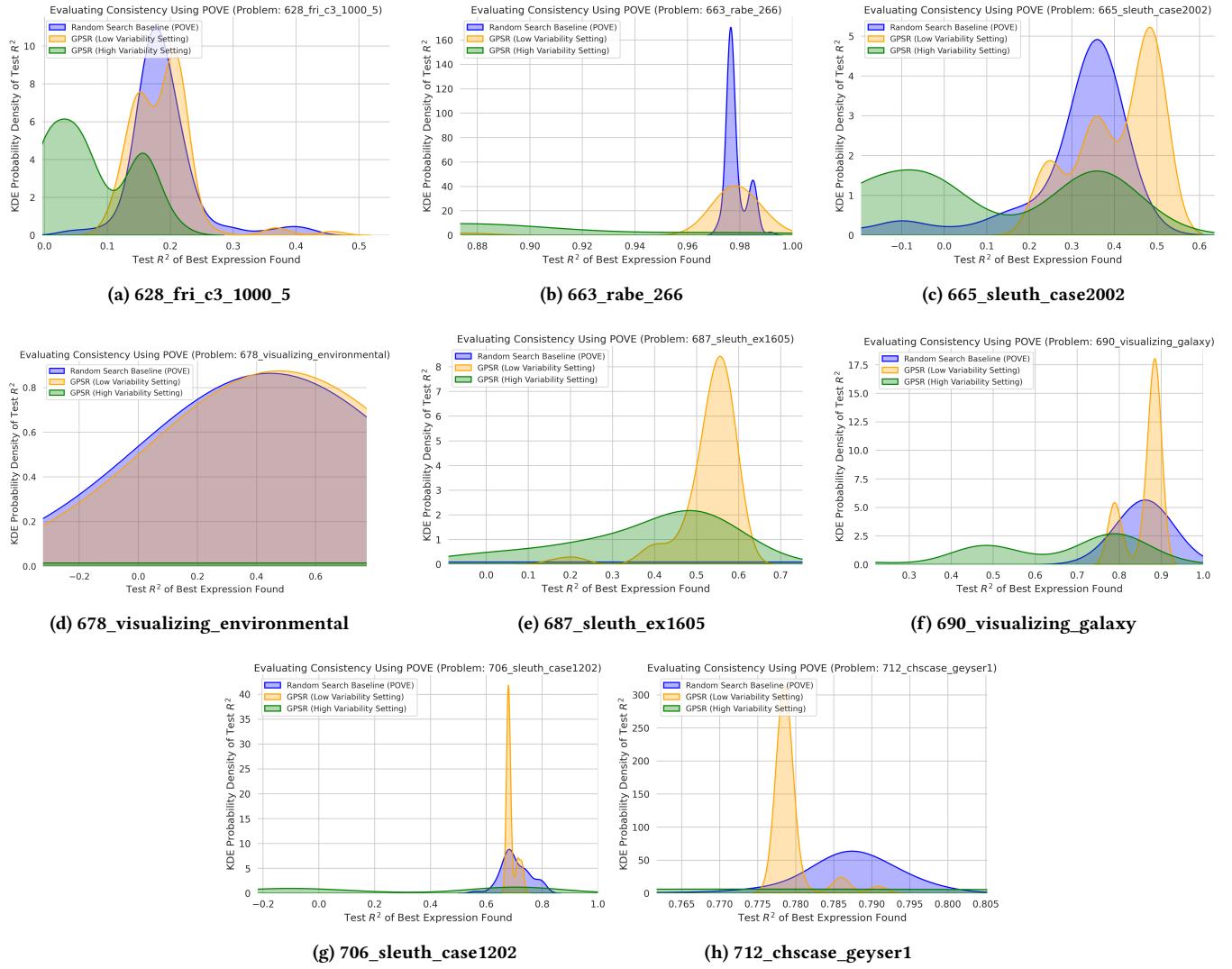


Figure 3: 25th to 32nd problem

B Other Benchmarking Plots

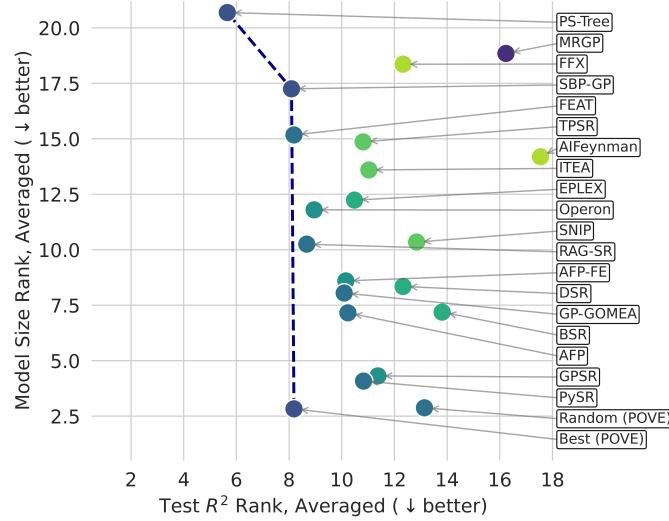


Figure 4: Rank-based Pareto frontier of mean model size and test R^2 . This is more common in SR literature, but ranking removes the difference in magnitude of the metrics. The Pareto frontier showed in the main paper is preferred.

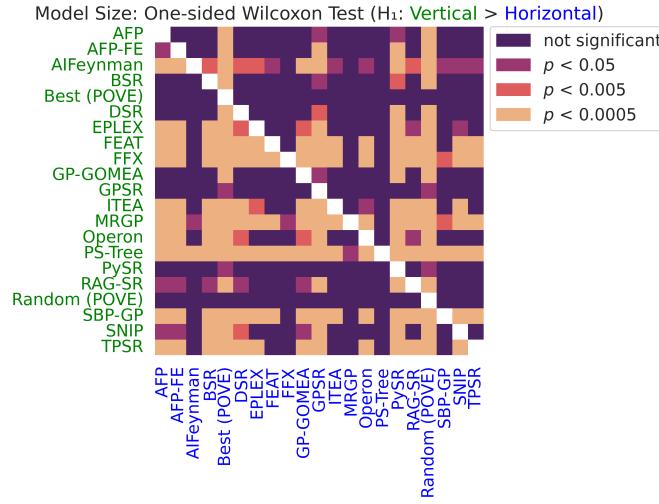


Figure 5: Pairwise one-sided Wilcoxon signed-rank test statistical test on model size.

C All Generalization and Optimal Parsimony Penalty Plots

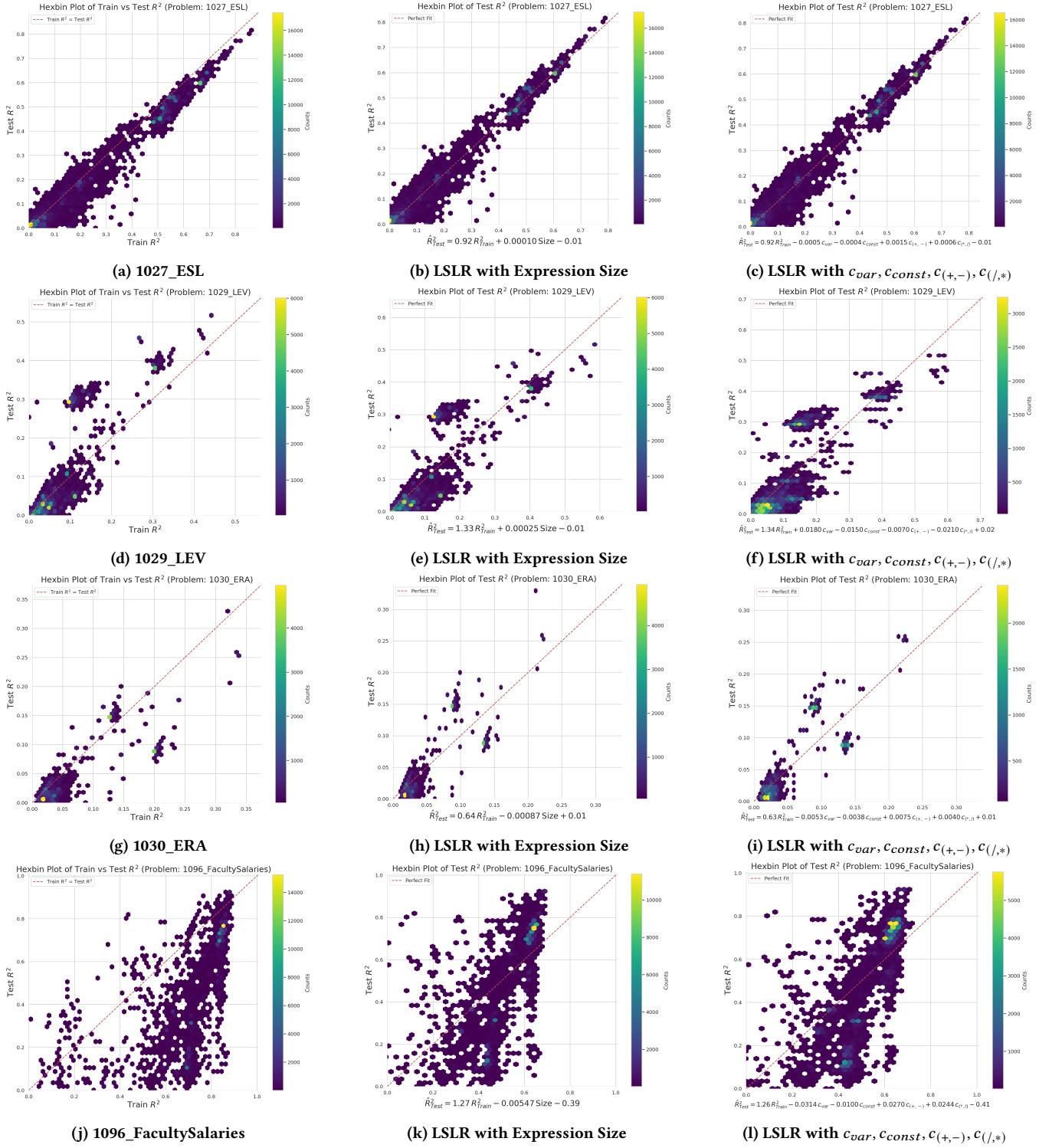


Figure 6: Problems from 1027_ESL to 1096_FacultySalaries



Figure 7: Problems from 192_vineyard to 230_machine_cpu

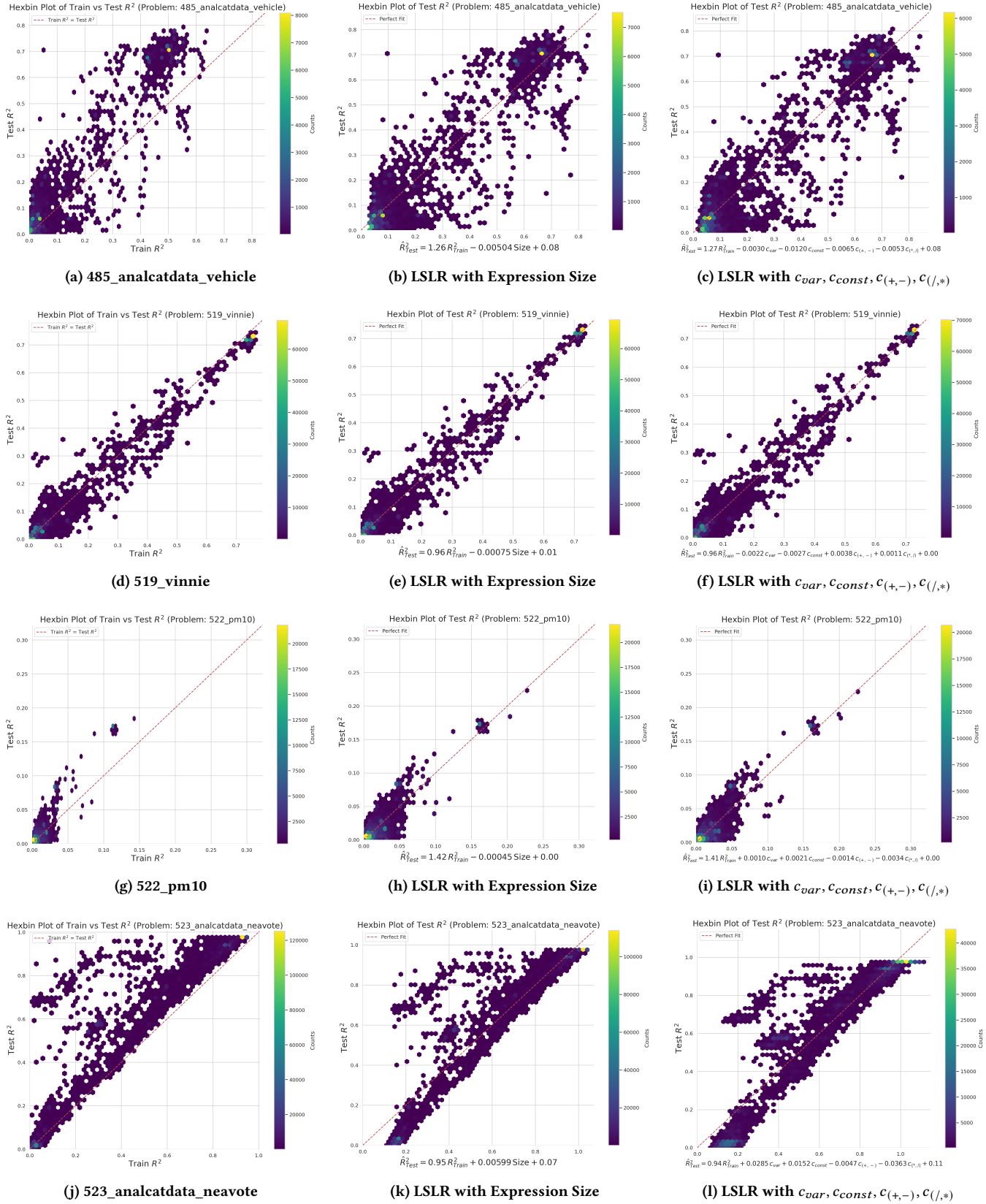


Figure 8: Problems from 485_analcatdata_vehicle to 523_analcatdata_neavote

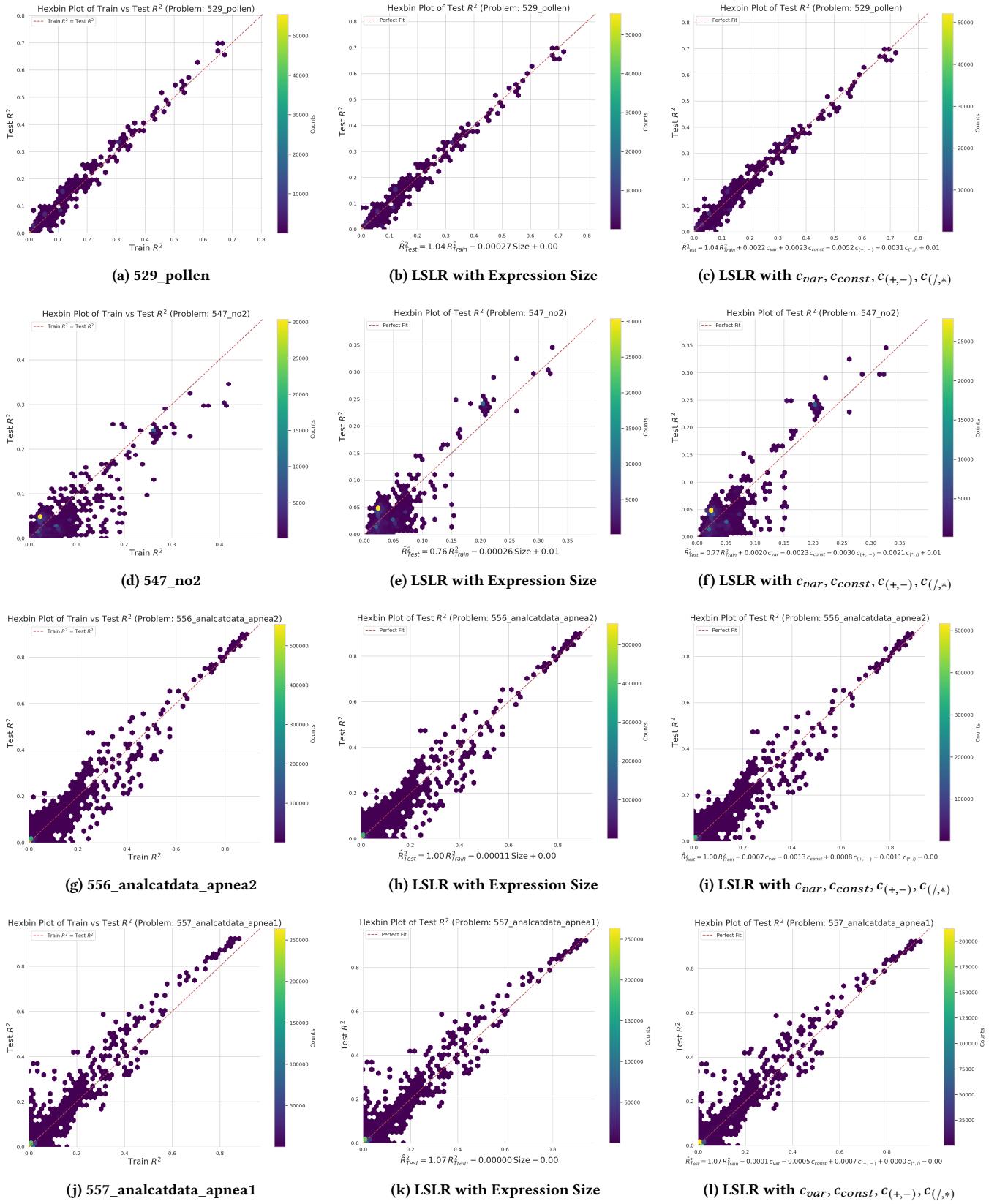


Figure 9: Problems from 529_pollen to 557_analcatdata_apnea1

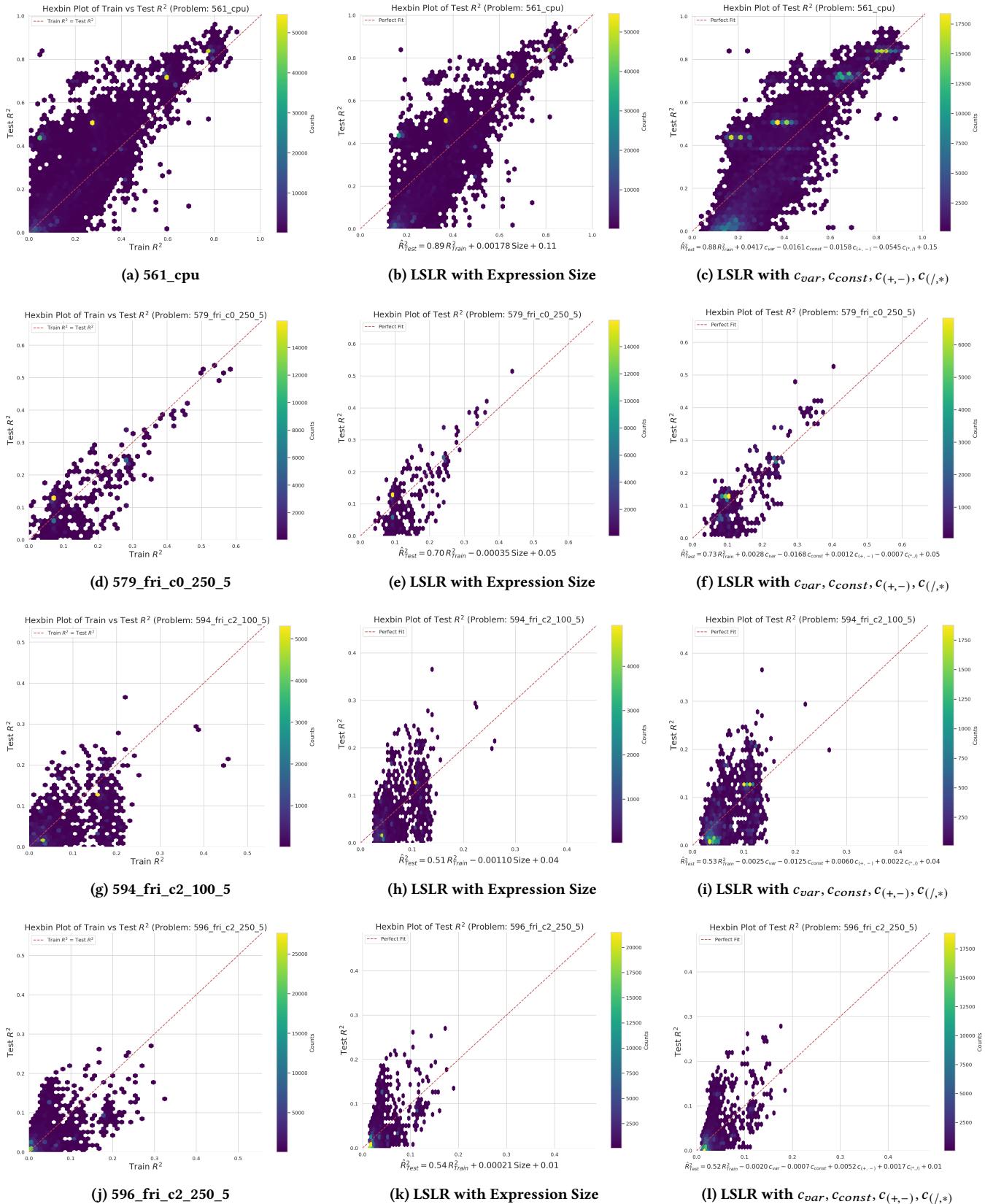


Figure 10: Problems from 561_cpu to 596_fri_c2_250_5

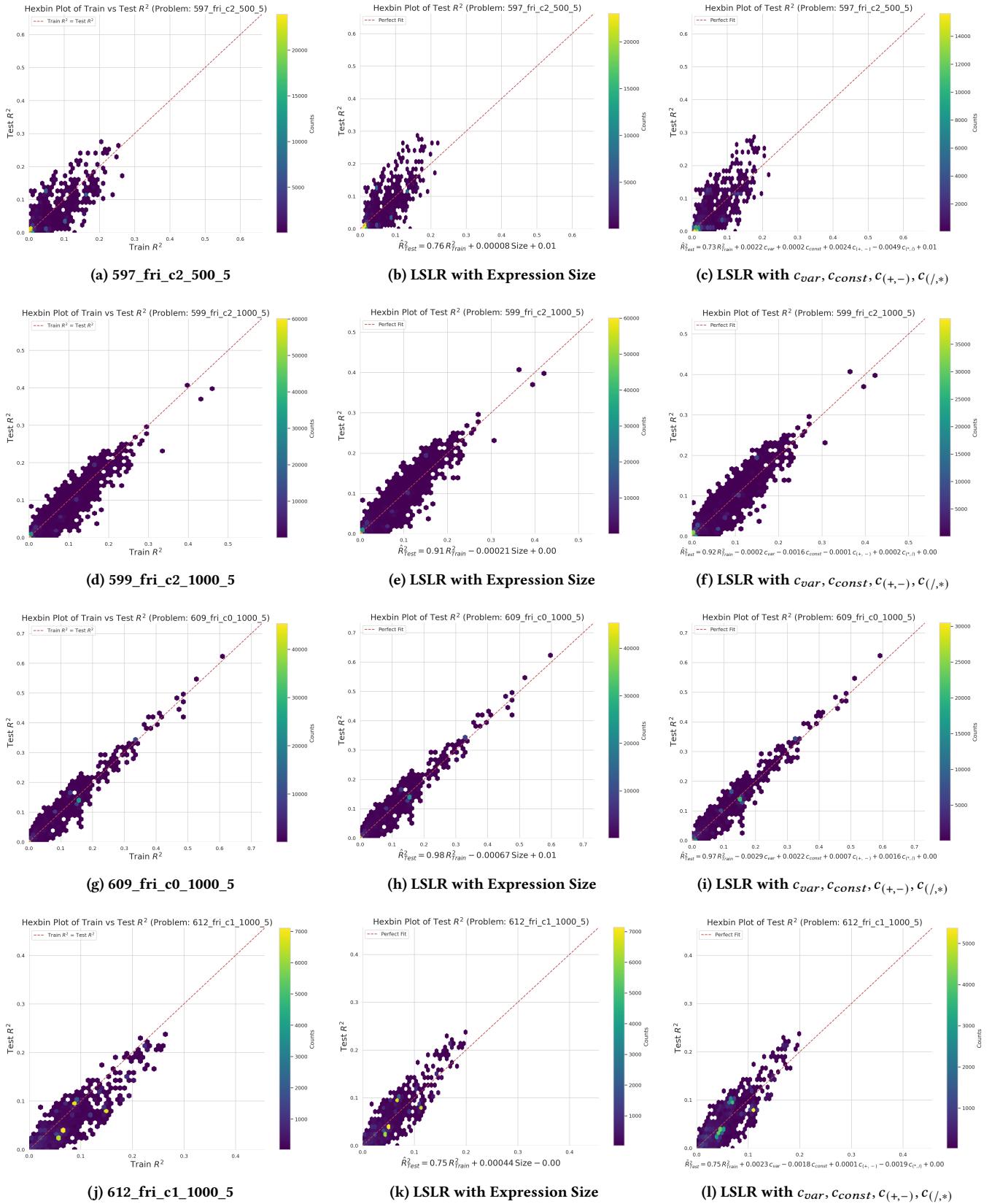


Figure 11: Problems from 597_fri_c2_500_5 to 612_fri_c1_1000_5

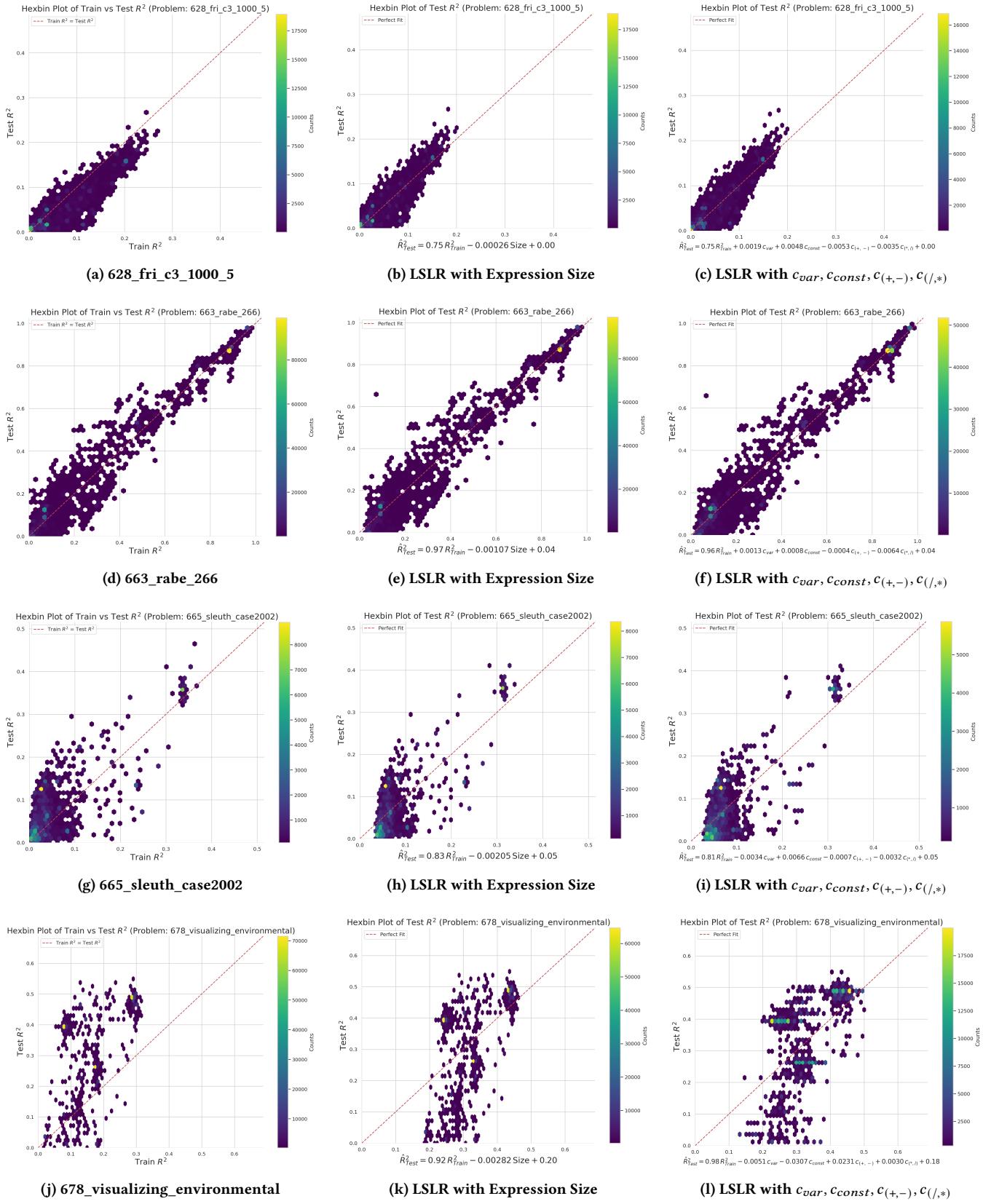


Figure 12: Problems from 628_fri_c3_1000_5 to 678_visualizing_environmental

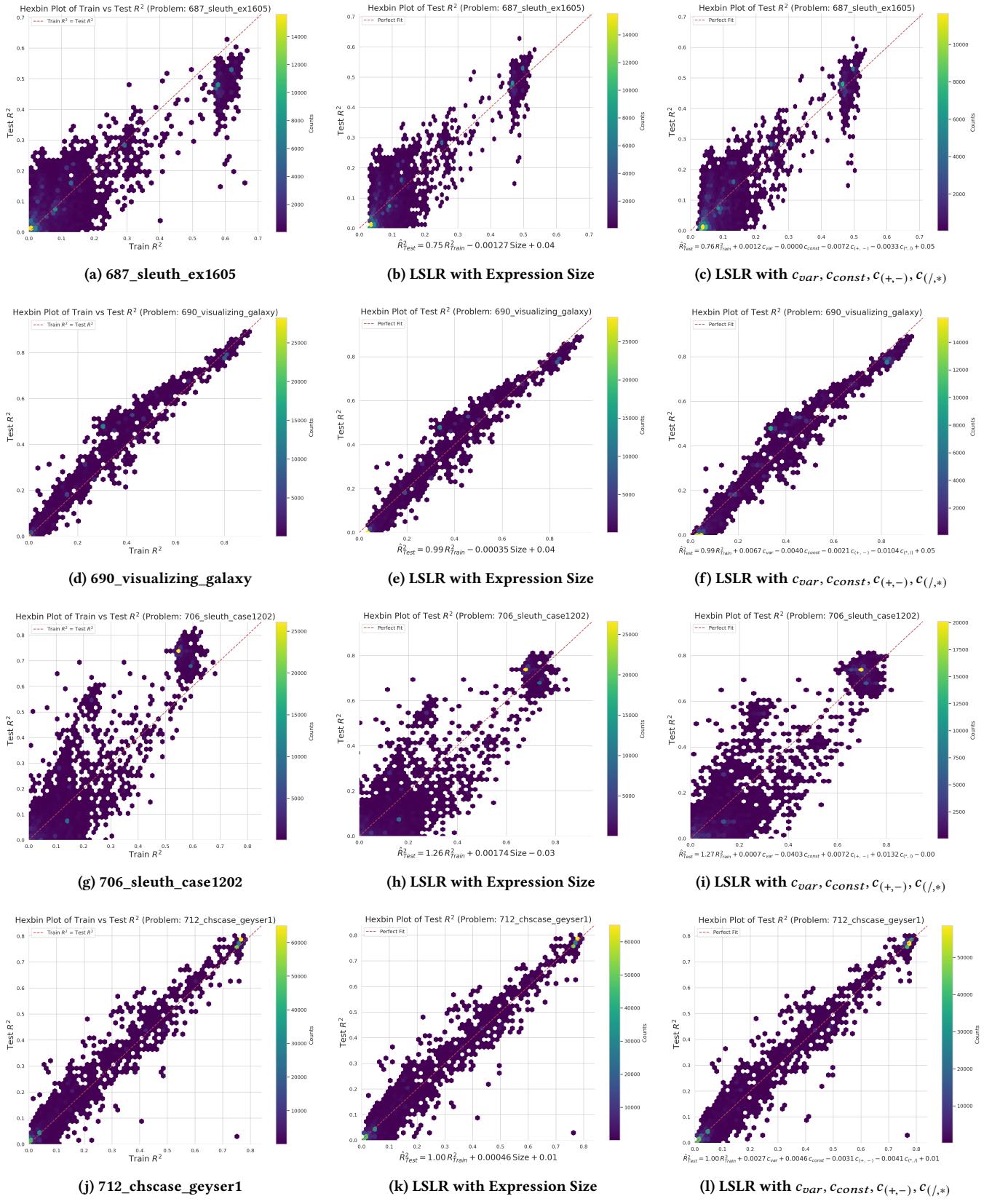


Figure 13: Problems from 687_sleuth_ex1605 to 712_chscase_geyser1

D All Distribution of Test R^2 Plots

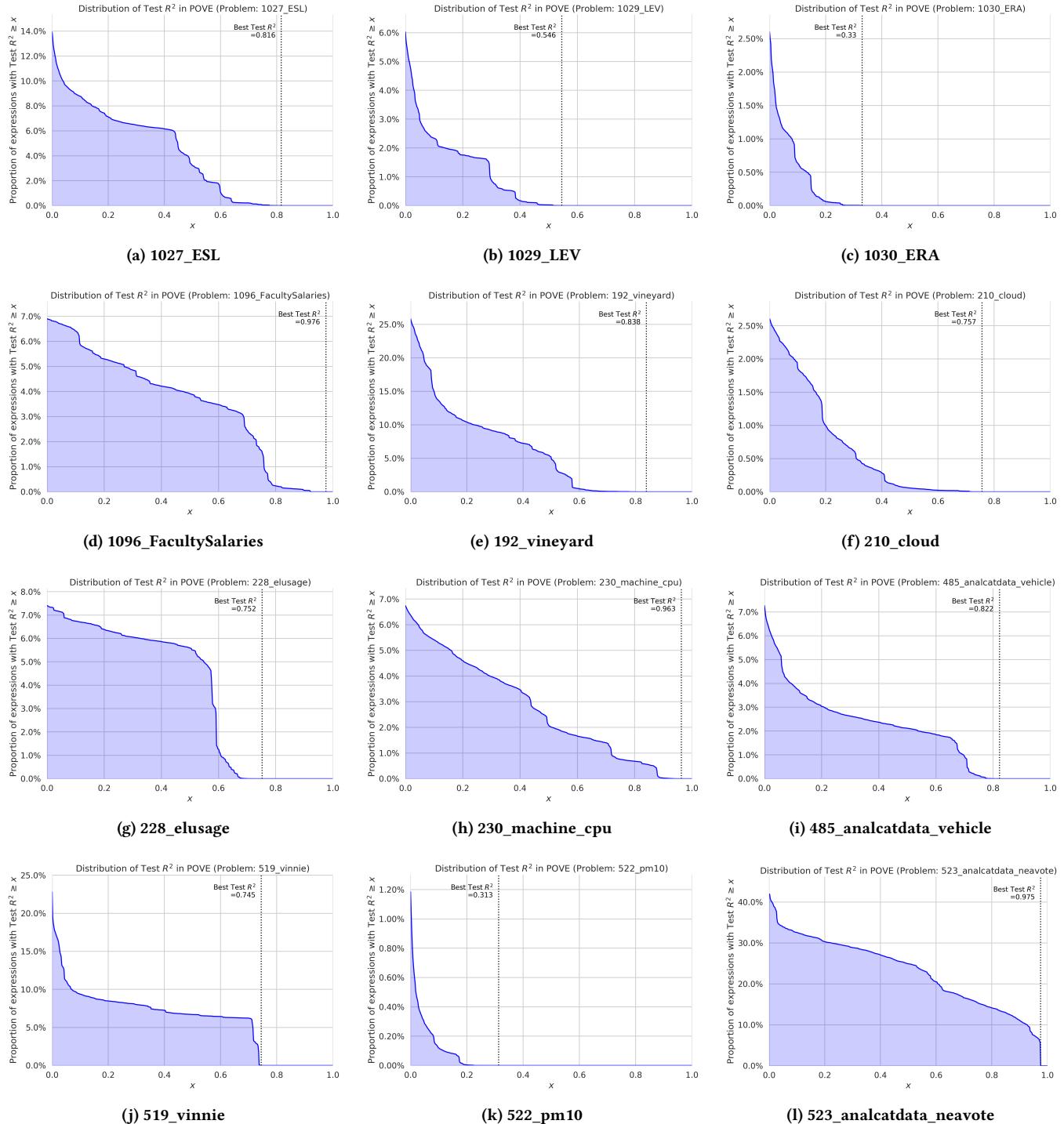


Figure 14: 1st to 12th problem

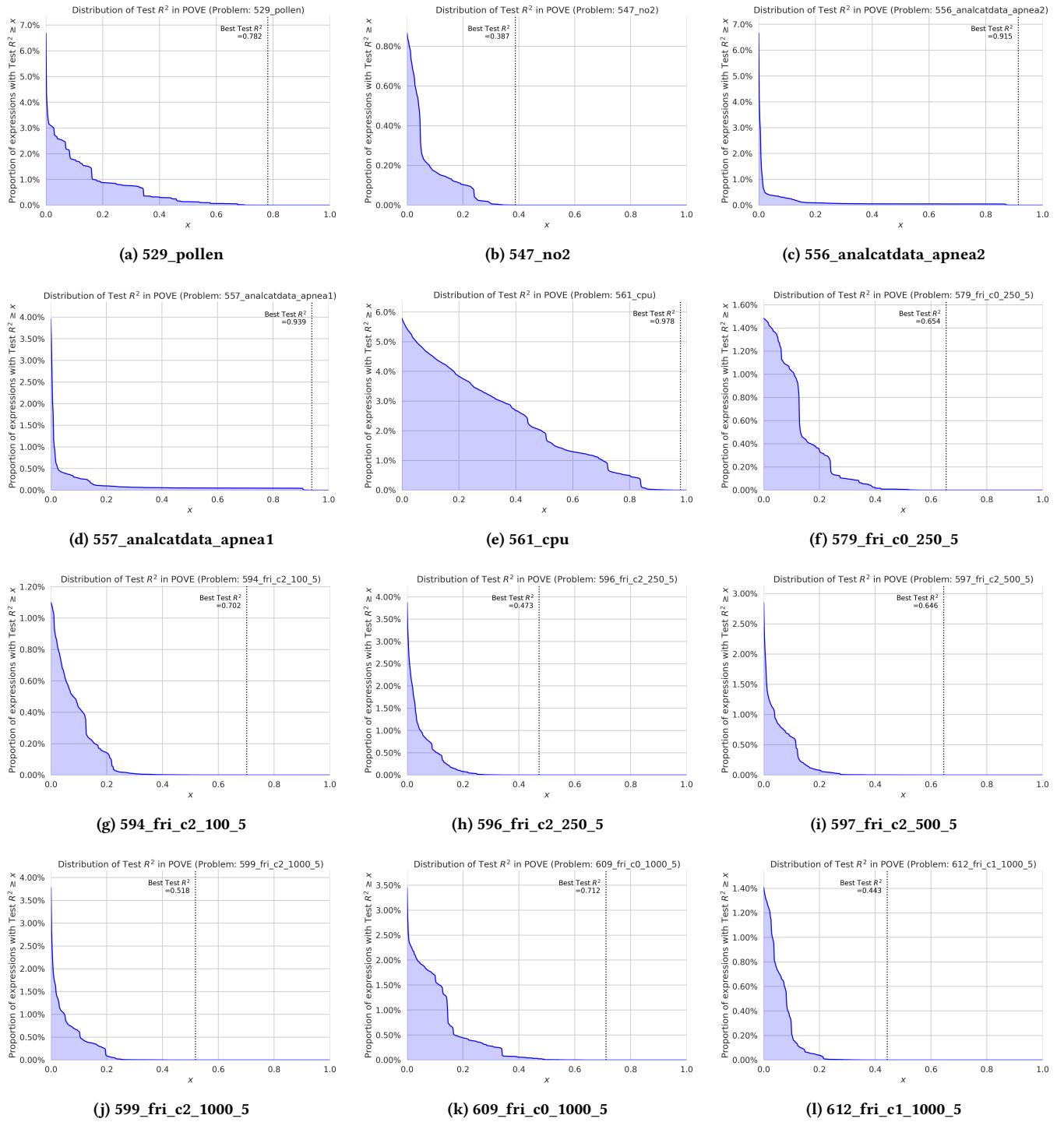


Figure 15: 13th to 24th problem

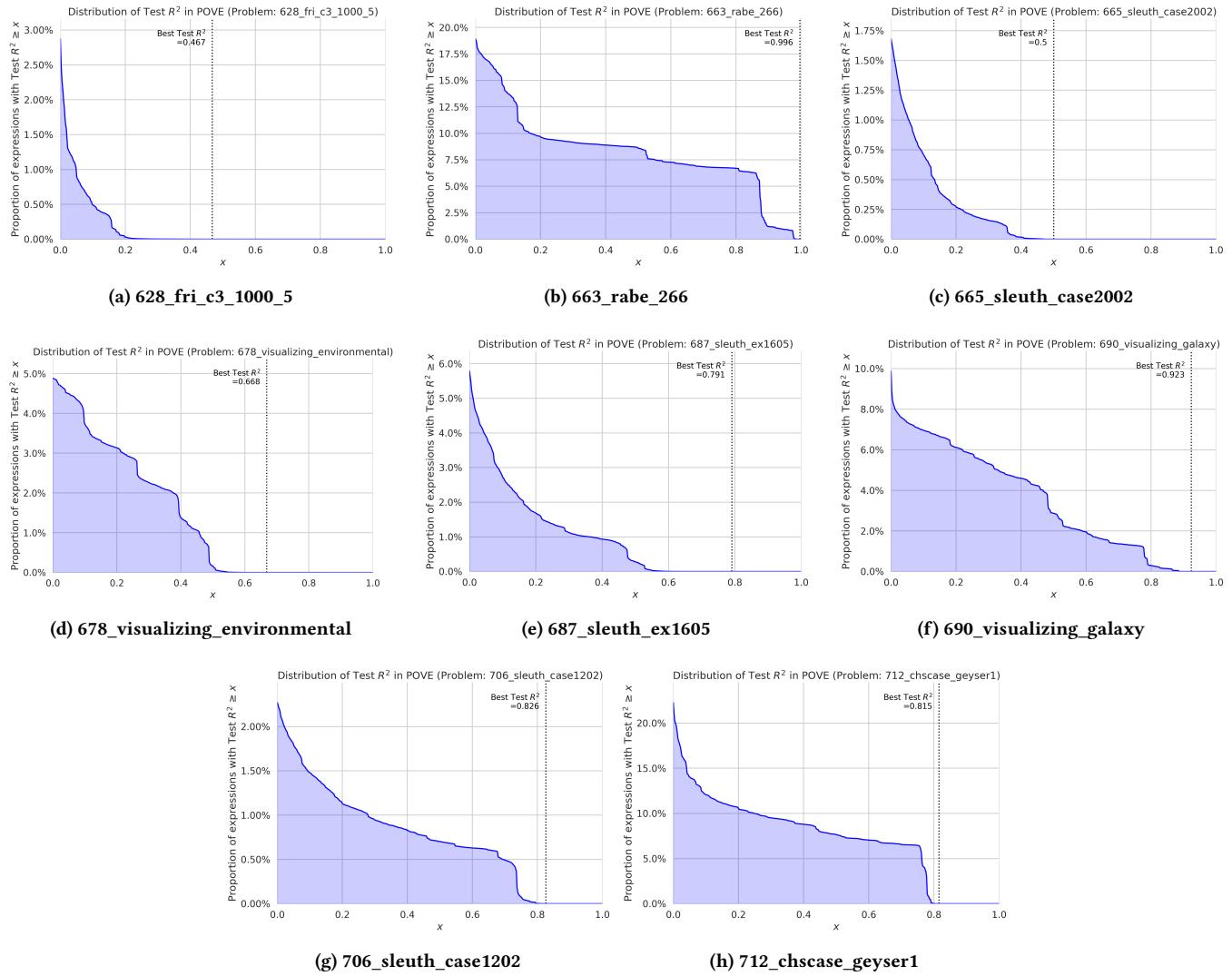


Figure 16: 25th to 32nd problem