

n	Name
1240	Complete dataset
748	Single mutation to alanine
273	Multiple mutations
130	Small-to-large mutation(s)
45	Multiple mutations, none alanine

Table 1: ZEMu dataset subset definition and composition

Mutation Category	Prediction Method	N	R	MAE	FC
Complete dataset	flex ddG	1240	<b>0.63</b>	<b>0.93</b>	<b>0.76</b>
	ddG monomer (hard-rep)		0.51	1.04	0.70
	no backrub control		0.56	1.00	0.74
	ZEMu paper		0.61	0.96	0.75
Small-to-large mutation(s)	flex ddG	130	<b>0.64</b>	<b>0.87</b>	<b>0.75</b>
	ddG monomer (hard-rep)		0.31	1.10	0.65
	no backrub control		0.41	1.00	0.72
	ZEMu paper		0.48	1.03	0.64
Mutation(s) to alanine	flex ddG	939	0.61	<b>0.89</b>	<b>0.77</b>
	ddG monomer (hard-rep)		0.50	0.98	0.71
	no backrub control		0.58	0.93	<b>0.77</b>
	ZEMu paper		<b>0.62</b>	0.90	<b>0.77</b>
Single mutation to alanine	flex ddG	748	<b>0.50</b>	<b>0.72</b>	<b>0.75</b>
	ddG monomer (hard-rep)		0.36	0.81	0.70
	no backrub control		0.44	0.78	<b>0.75</b>
	ZEMu paper		0.45	0.76	<b>0.75</b>
Multiple mutations	flex ddG	273	0.62	1.51	0.77
	ddG monomer (hard-rep)		0.50	1.69	0.66
	no backrub control		0.58	1.59	0.72
	ZEMu paper		<b>0.64</b>	<b>1.46</b>	<b>0.78</b>
Multiple mutations, all alanine	flex ddG	191	0.47	1.55	<b>0.85</b>
	ddG monomer (hard-rep)		0.34	1.67	0.74
	no backrub control		0.50	1.51	0.83
	ZEMu paper		<b>0.55</b>	<b>1.44</b>	<b>0.85</b>
Multiple mutations, none alanine	flex ddG	45	<b>0.67</b>	<b>1.57</b>	<b>0.53</b>
	ddG monomer (hard-rep)		0.40	1.96	0.49
	no backrub control		0.44	1.82	<b>0.53</b>
	ZEMu paper		0.53	1.79	0.51
Antibodies	flex ddG	355	<b>0.60</b>	<b>0.89</b>	<b>0.74</b>
	ddG monomer (hard-rep)		0.50	0.98	0.71
	no backrub control		0.49	0.96	0.73
	ZEMu paper		0.54	0.96	<b>0.74</b>

Table 2: Main results table. Backrub steps = 35000. R = Pearson’s R. MAE = Mean Absolute Error. FC = Fraction Correct. N = number of mutations in the dataset or subset.

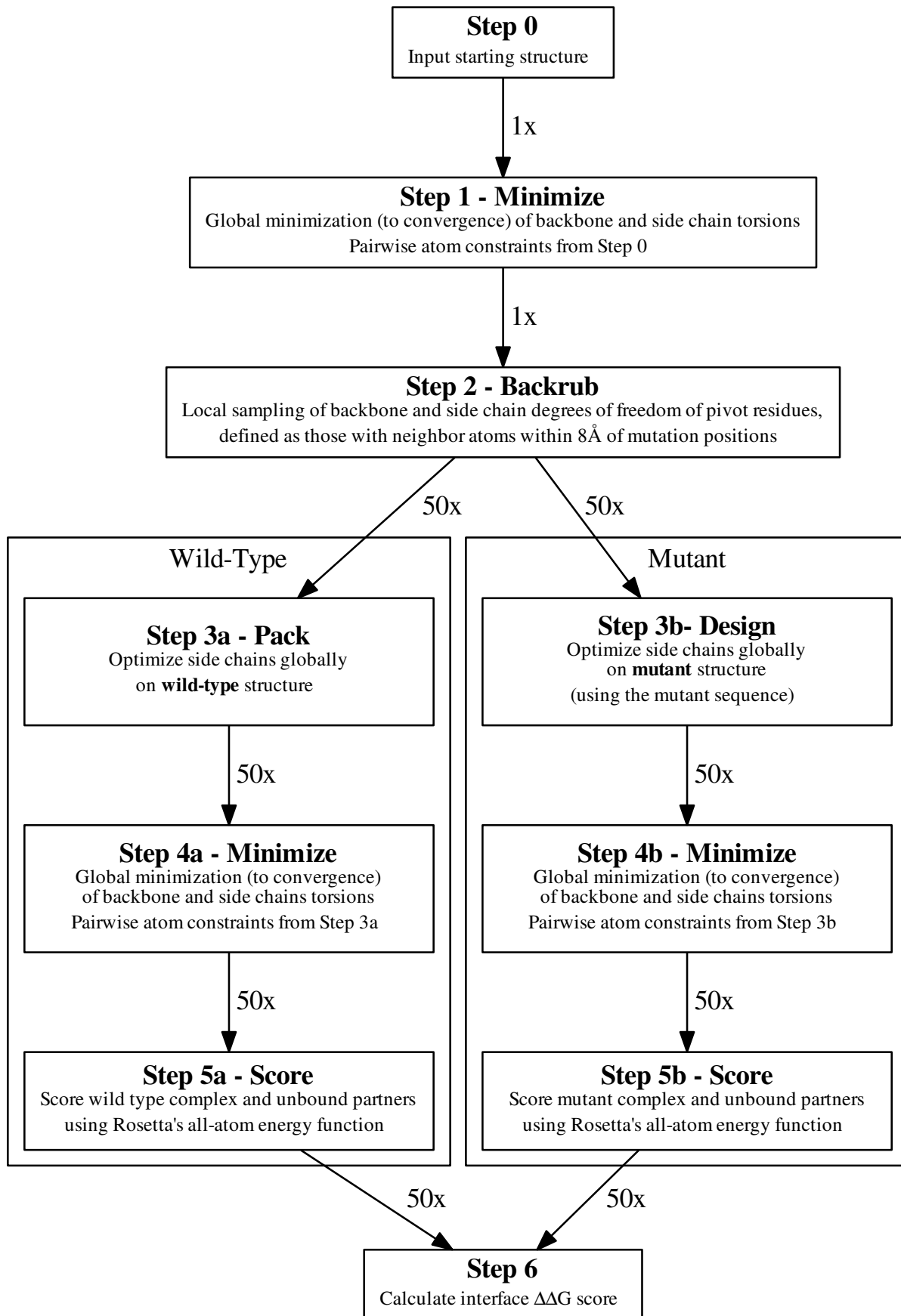


Figure 1: Schematic of the Alex ddG protocol method.

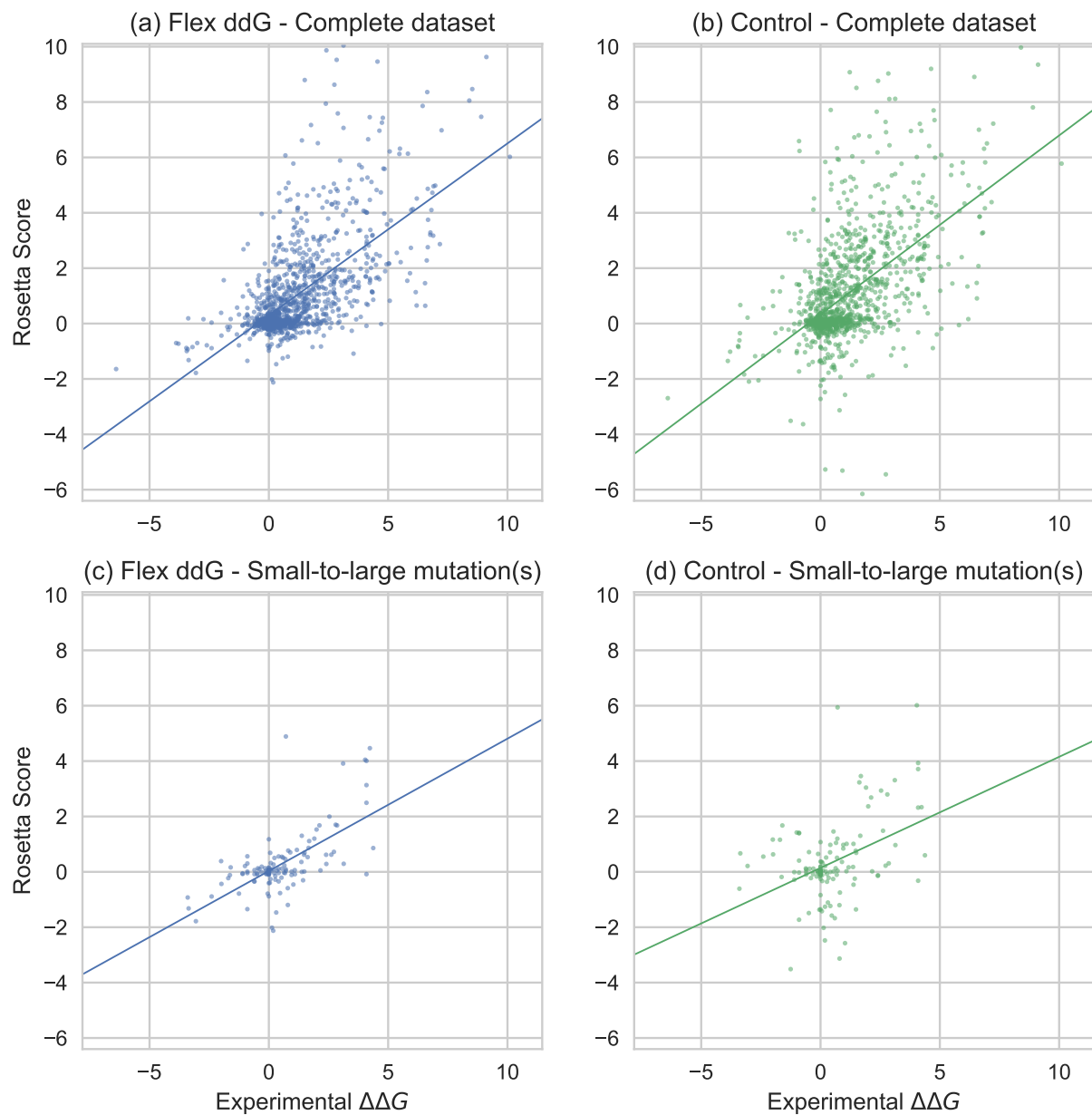


Figure 2: Experimentally determined  $\Delta\Delta G$  values (y-axis) vs. Rosetta predictions. (a) Flex ddg method (35000 backrub steps); Complete dataset mutation set (n=1240). (b) No backrub control; Complete dataset mutation set (n=1240). (c) Flex ddg method (35000 backrub steps); Small-to-large mutation(s) mutation set (n=130). (d) No backrub control; Small-to-large mutation(s) mutation set (n=130).

# $\Delta\Delta G$ prediction performance vs. number of backrub sampling steps

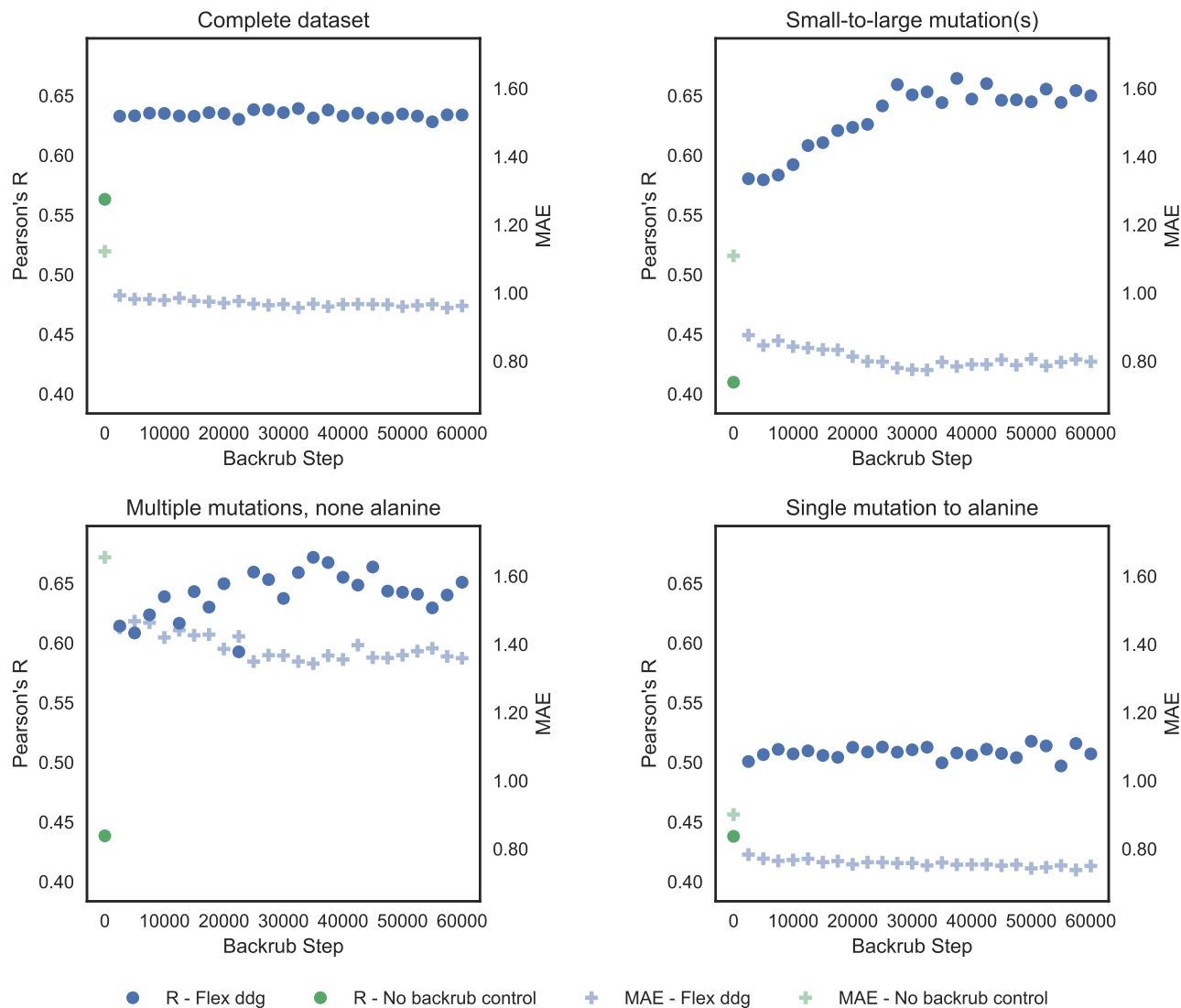


Figure 3: Correlation (Pearson's R) and MAE (Mean Absolute Error) vs. number of backrub steps, on the complete ZEMu set, and subsets. (a) Complete dataset (n=1240) (b) Small-to-large mutation(s) (n=130) (c) Multiple mutations, none alanine (n=45) (d) Single mutation to alanine (n=748)

Mutation Category	Prediction Method	N	R	MAE	FC
Complete dataset	flex ddG	1240	<b>0.64</b>	0.93	<b>0.76</b>
	flex ddG (1.6 kT)		<b>0.64</b>	<b>0.92</b>	<b>0.76</b>
Small-to-large mutation(s)	flex ddG	130	0.59	0.88	0.72
	flex ddG (1.6 kT)		<b>0.64</b>	<b>0.85</b>	<b>0.74</b>
Single mutation to alanine	flex ddG	748	<b>0.51</b>	0.73	<b>0.75</b>
	flex ddG (1.6 kT)		<b>0.51</b>	<b>0.72</b>	<b>0.75</b>
Multiple mutations	flex ddG	273	<b>0.63</b>	1.51	<b>0.79</b>
	flex ddG (1.6 kT)		<b>0.63</b>	<b>1.50</b>	<b>0.79</b>
Multiple mutations, none alanine	flex ddG	45	<b>0.64</b>	1.62	<b>0.56</b>
	flex ddG (1.6 kT)		0.62	<b>1.61</b>	0.53

Table 3: Flex ddG performance comparison, when backrub is run with a sampling temperature (kT) of 1.2 or 1.6. Backrub steps = 10000. R = Pearson's R. MAE = Mean Absolute Error. FC = Fraction Correct. N = number of mutations in the dataset or subset.

# $\Delta\Delta G$ prediction performance vs. number of backrub sampling steps

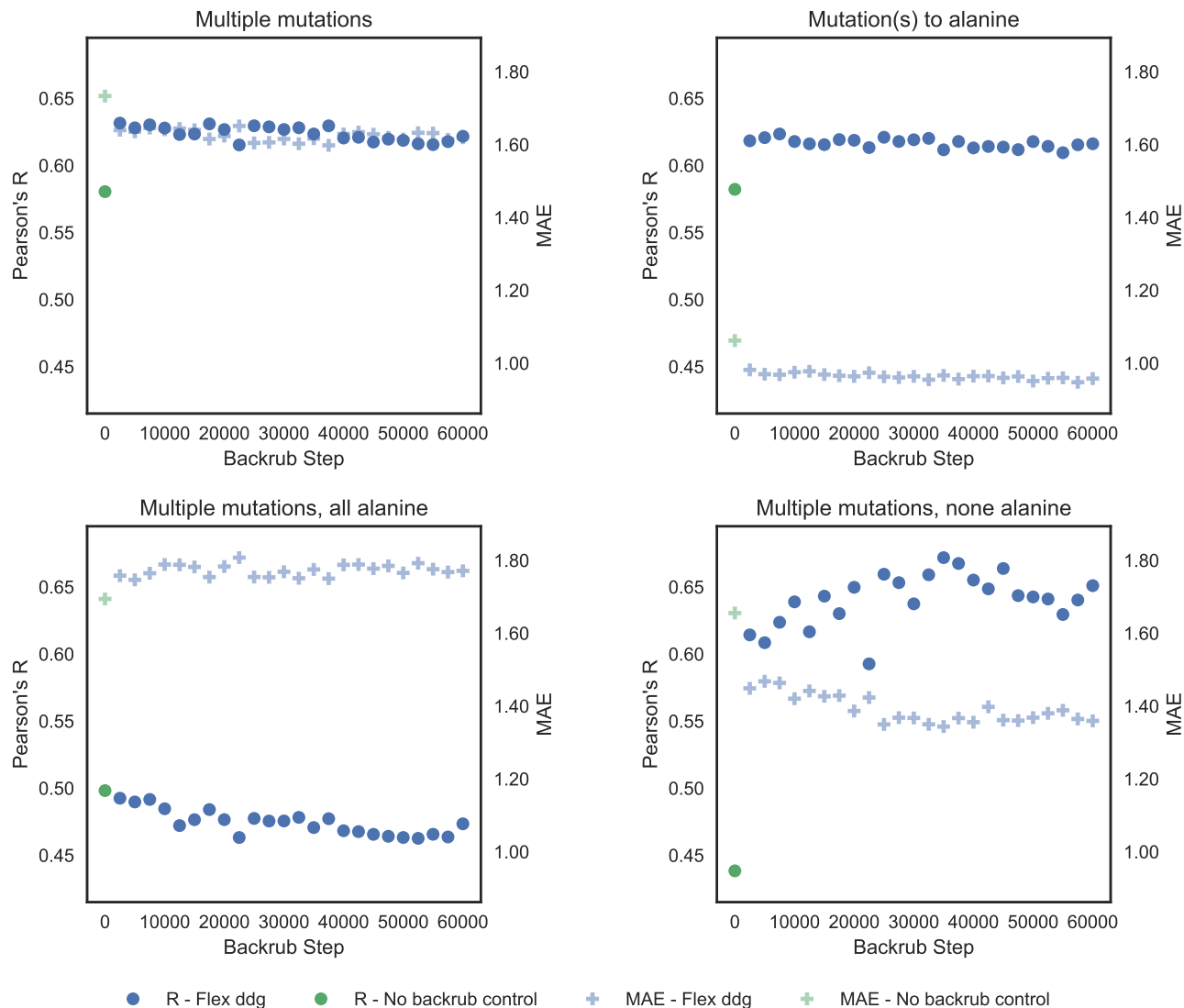


Figure 4: Correlation (Pearson’s R) and MAE (Mean Absolute Error) vs. number of backrub steps, on the complete ZEMu set, and subsets. (a) Multiple mutations (n=273) (b) Mutation(s) to alanine (n=939) (c) Multiple mutations, all alanine (n=191) (d) Multiple mutations, none alanine (n=45)

Mutation Category	Prediction Method	N	R	MAE	FC
Multiple mutations	flex ddG	273	0.62	1.51	0.77
	no backrub control		0.58	1.59	0.72
	ZEMu paper		<b>0.64</b>	<b>1.46</b>	<b>0.78</b>
Multiple mutations, all alanine	flex ddG	191	0.47	1.55	<b>0.85</b>
	no backrub control		0.50	1.51	0.83
	ZEMu paper		<b>0.55</b>	<b>1.44</b>	<b>0.85</b>
Multiple mutations, none alanine	flex ddG	45	<b>0.67</b>	<b>1.57</b>	<b>0.53</b>
	no backrub control		0.44	1.82	<b>0.53</b>
	ZEMu paper		0.53	1.79	0.51
Mutation(s) to alanine	flex ddG	939	0.61	<b>0.89</b>	<b>0.77</b>
	no backrub control		0.58	0.93	<b>0.77</b>
	ZEMu paper		<b>0.62</b>	0.90	<b>0.77</b>

Table 4: Multiple mutations results (backrub steps = 35000). R = Pearson’s R. MAE = Mean Absolute Error. FC = Fraction Correct. N = number of mutations in the dataset or subset.

# $\Delta\Delta G$ prediction performance vs. number of backrub sampling steps

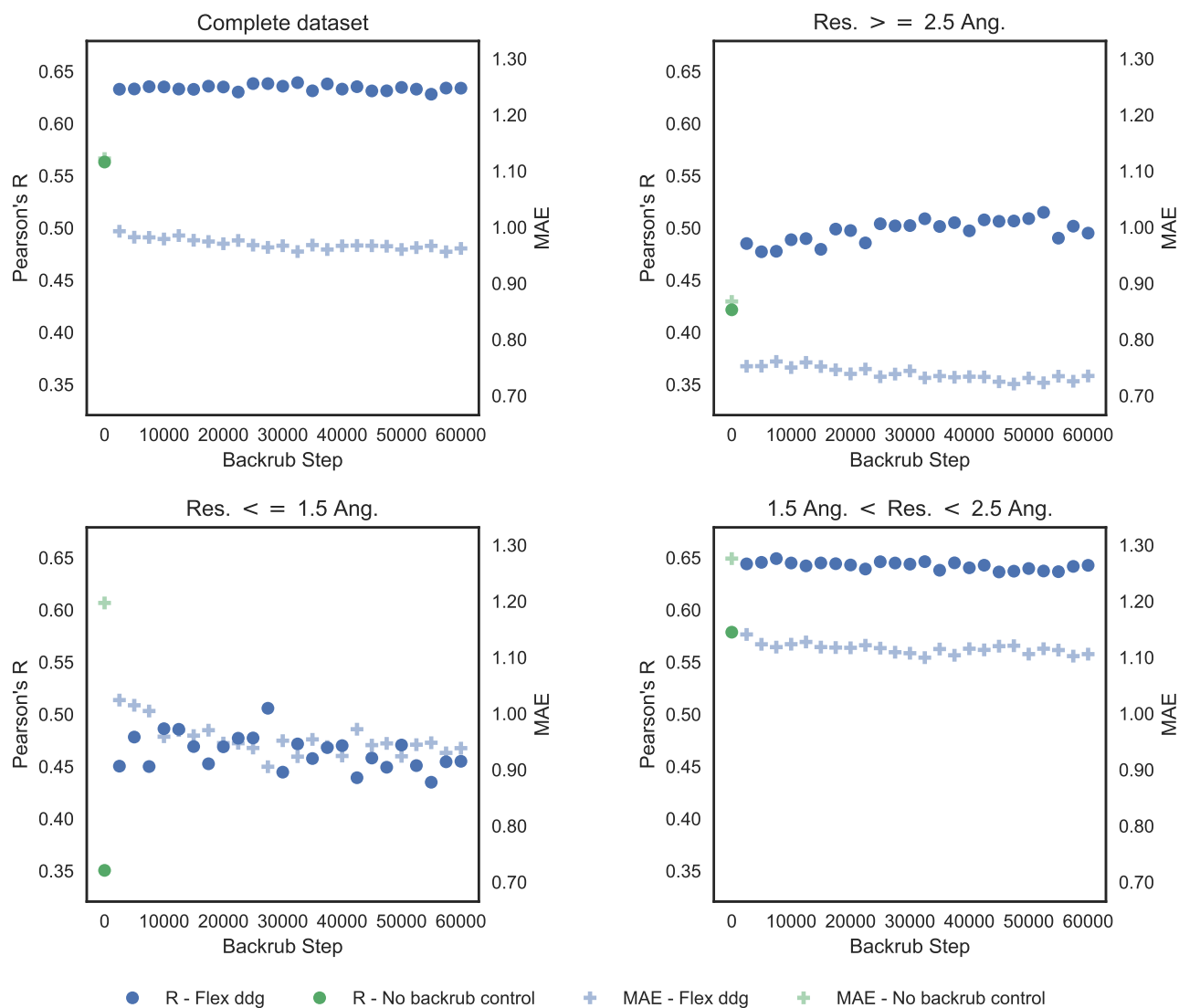


Figure 5: Correlation (Pearson’s R) and MAE (Mean Absolute Error) vs. number of backrub steps, on the complete ZEMu set, and subsets. (a) Complete dataset (n=1240) (b) Res.  $\geq 2.5$  ang. (n=457) (c) Res.  $\leq 1.5$  ang. (n=52) (d) 1.5 ang.  $<$  res.  $< 2.5$  ang. (n=731)

Mutation Category	Prediction Method	N	R	MAE	FC
Complete dataset	flex ddG	1240	<b>0.63</b>	<b>0.93</b>	0.76
	no backrub control		0.56	1.00	0.74
	ddG monomer (hard-rep)		0.62	0.94	<b>0.77</b>
	ZEMu paper		0.61	0.96	0.75
Antibodies	flex ddG	355	<b>0.60</b>	<b>0.89</b>	0.74
	no backrub control		0.49	0.96	0.73
	ddG monomer (hard-rep)		0.58	0.90	<b>0.77</b>
	ZEMu paper		0.54	0.96	0.74

Table 5: Performance of the Rosetta flex ddG method on the subset of complexes containing an antibody binding partner (backrub steps = 35000). R = Pearson’s R. MAE = Mean Absolute Error. FC = Fraction Correct. N = number of mutations in the dataset or subset.

# $\Delta\Delta G$ prediction performance vs. number of backrub sampling steps

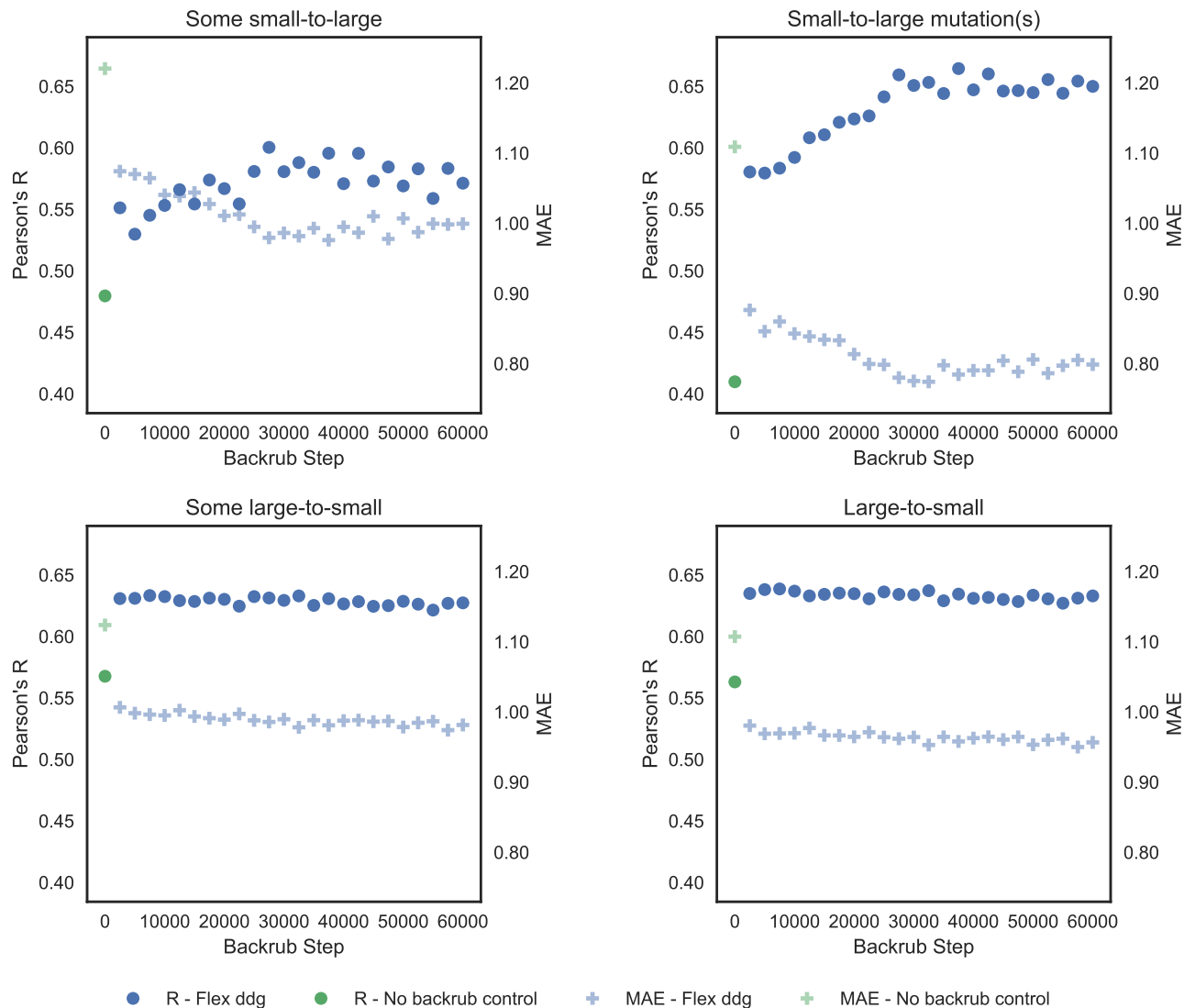


Figure 6: Correlation (Pearson's R) and MAE (Mean Absolute Error) vs. number of backrub steps, on the complete ZEMu set, and subsets. (a) Some small-to-large (n=164) (b) Small-to-large mutation(s) (n=130) (c) Some large-to-small (n=1110) (d) Large-to-small (n=1076)

Git SHA1	Protocol
69aa5266f0d5	flex ddG
69aa5266f0d5	no backrub control
3b2aa5cc3798	ddG monomer

Table 6: SHA1 Git version of Rosetta used for benchmarking

## $\Delta\Delta G$ prediction performance vs. number of structural ensemble members

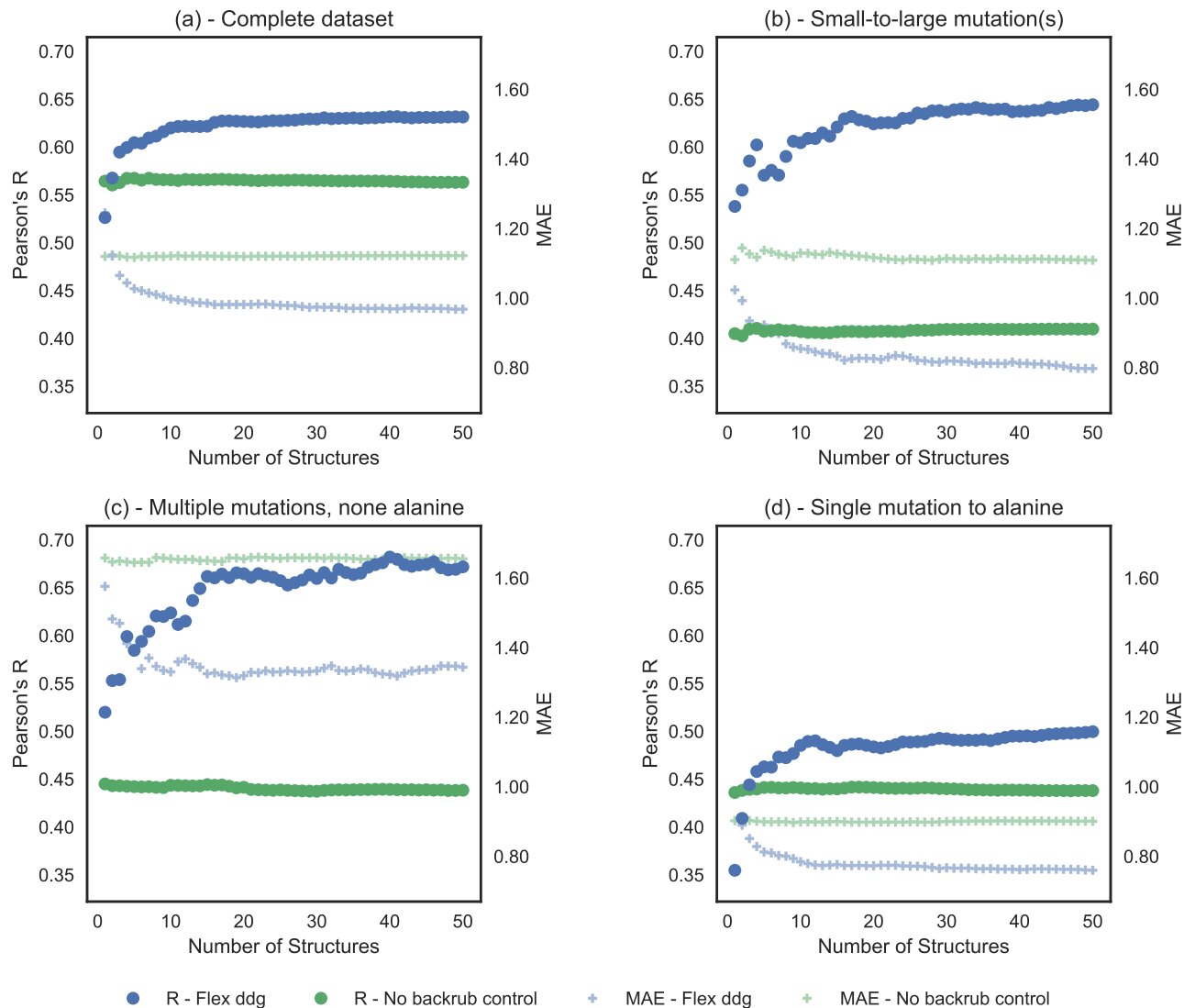


Figure 7: Correlation (Pearson's R) and MAE (Mean Absolute Error) vs. number of averaged structures, on the complete ZEMu set, and subsets. Structures are not sorted, and are randomly added to the ensemble. (a) Complete dataset ( $n = 1240$ , backrub steps = 35000) (b) Small-to-large mutation(s) ( $n = 130$ , backrub steps = 35000) (c) Multiple mutations, none alanine ( $n = 45$ , backrub steps = 35000) (d) Single mutation to alanine ( $n = 748$ , backrub steps = 35000)



## $\Delta\Delta G$ prediction performance vs. number of structural ensemble members

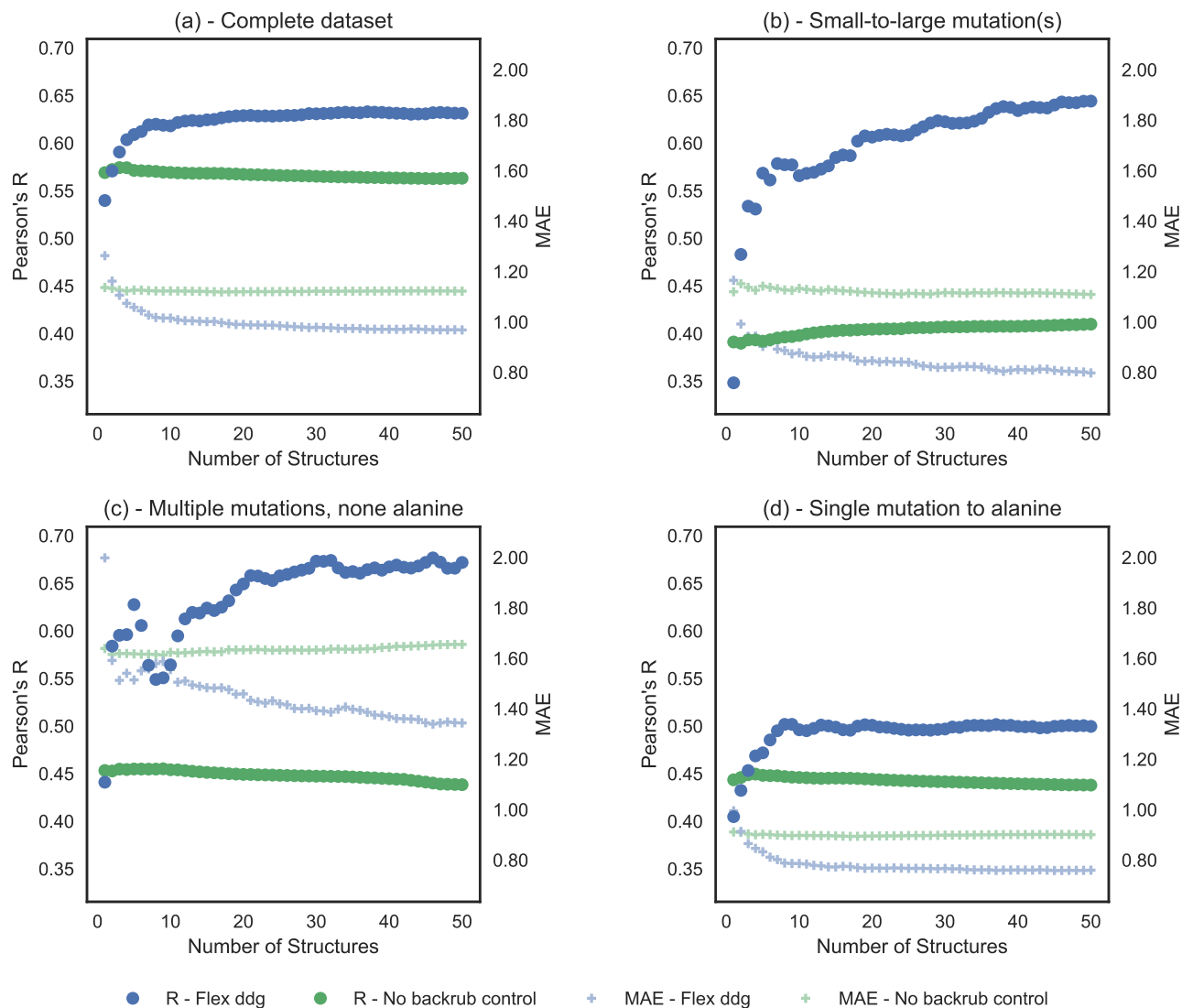


Figure 8: Correlation (Pearson's R) and MAE (Mean Absolute Error) vs. number of averaged structures, on the complete ZEMu set, and subsets. Structures are sorted by their minimized wild-type complex energy. (a) Complete dataset ( $n = 1240$ , backrub steps = 35000) (b) Small-to-large mutation(s) ( $n = 130$ , backrub steps = 35000) (c) Multiple mutations, none alanine ( $n = 45$ , backrub steps = 35000) (d) Single mutation to alanine ( $n = 748$ , backrub steps = 35000)

## $\Delta\Delta G$ prediction performance vs. number of structural ensemble members

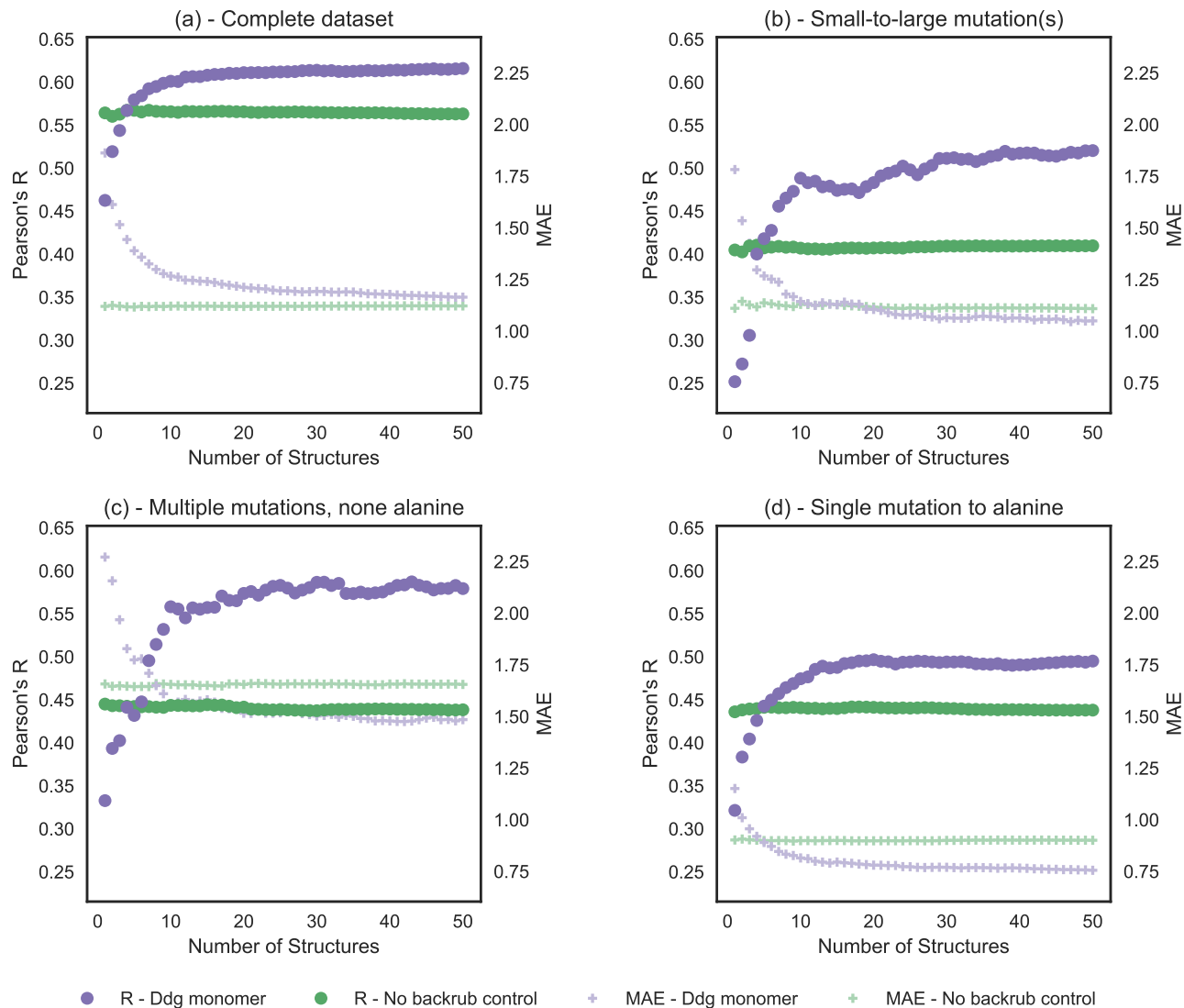


Figure 9: Correlation (Pearson's R) and MAE (Mean Absolute Error) vs. number of averaged structures, on the complete ZEMu set, and subsets. Structures are not sorted, and are randomly added to the ensemble. (a) Complete dataset (n = 1240) (b) Small-to-large mutation(s) (n = 130) (c) Multiple mutations, none alanine (n = 45) (d) Single mutation to alanine (n = 748)

## $\Delta\Delta G$ prediction performance vs. number of structural ensemble members

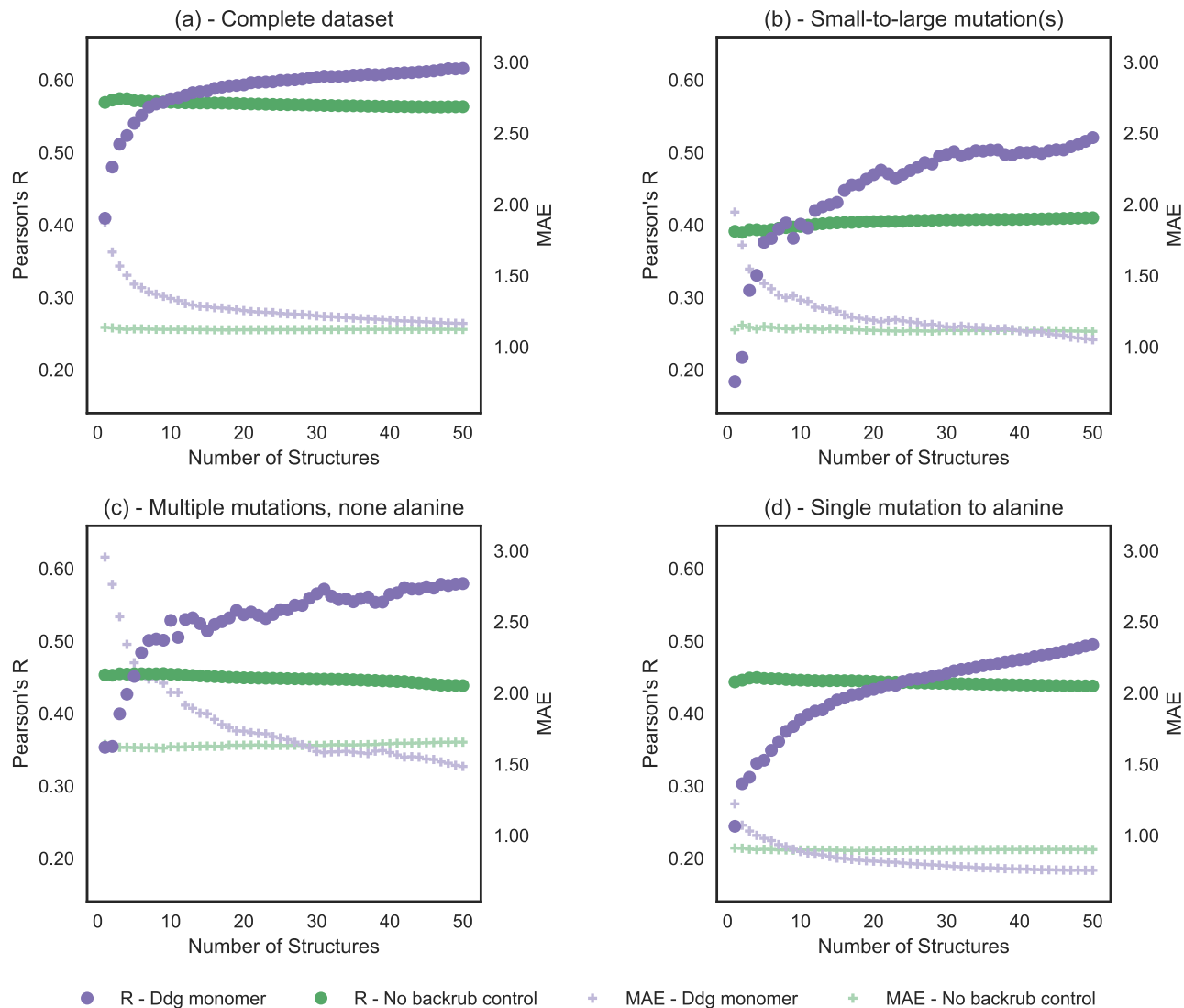


Figure 10: Correlation (Pearson's R) and MAE (Mean Absolute Error) vs. number of averaged structures, on the complete ZEMu set, and subsets. Structures are sorted by their minimized wild-type complex energy. (a) Complete dataset (n = 1240) (b) Small-to-large mutation(s) (n = 130) (c) Multiple mutations, none alanine (n = 45) (d) Single mutation to alanine (n = 748)

Mutation Category	Prediction Method	N	R	MAE	FC
Stabilizing	flex ddG	32	<b>0.62</b>	2.65	0.00
	no backrub control		0.50	2.83	<b>0.03</b>
	ddG monomer (hard-rep)		0.39	2.66	<b>0.03</b>
	ZEMu paper		0.31	<b>2.64</b>	<b>0.03</b>
Neutral	flex ddG	719	<b>0.19</b>	<b>0.58</b>	<b>0.85</b>
	no backrub control		0.10	0.68	0.81
	ddG monomer (hard-rep)		0.13	0.61	0.84
	ZEMu paper		0.16	0.64	0.81
Positive	flex ddG	489	<b>0.48</b>	1.32	0.67
	no backrub control		0.44	1.35	0.70
	ddG monomer (hard-rep)		0.47	<b>1.31</b>	<b>0.71</b>
	ZEMu paper		<b>0.48</b>	1.32	0.70

Table 7: Performance of the Rosetta flex ddG method on the subset of mutations experimentally determined to be stabilizing ( $\Delta\Delta G < 0$ ). Backrub steps = 35000. R = Pearson’s R. MAE = Mean Absolute Error. FC = Fraction Correct. N = number of mutations in the dataset or subset.

Mutation Category	Prediction Method	N	R	MAE	FC
pdb-1A22	flex ddG	142	<b>0.33</b>	<b>0.63</b>	0.75
	no backrub control		0.18	0.72	0.74
	ddG monomer (hard-rep)		0.12	0.70	0.74
	ZEMu paper		0.19	0.64	<b>0.79</b>
pdb-1A4Y	flex ddG	45	0.80	1.46	0.76
	no backrub control		0.79	1.64	<b>0.80</b>
	ddG monomer (hard-rep)		0.77	1.55	0.69
	ZEMu paper		<b>0.87</b>	<b>1.38</b>	0.71
pdb-1ACB	flex ddG	6	0.25	3.42	0.83
	no backrub control		0.23	3.59	<b>1.00</b>
	ddG monomer (hard-rep)		0.58	3.27	<b>1.00</b>
	ZEMu paper		<b>0.79</b>	<b>2.82</b>	<b>1.00</b>
pdb-1AHW	flex ddG	10	-0.77	1.12	0.4
	no backrub control		-0.42	1.15	0.4
	ddG monomer (hard-rep)		-0.34	1.03	0.4
	ZEMu paper		<b>0.30</b>	<b>0.92</b>	<b>0.6</b>
pdb-1AK4	flex ddG	15	<b>0.69</b>	<b>0.53</b>	0.73
	no backrub control		0.35	0.68	0.67
	ddG monomer (hard-rep)		0.63	0.69	<b>0.80</b>
	ZEMu paper		0.44	1.04	0.53
pdb-1CBW	flex ddG	15	-0.08	0.55	0.60
	no backrub control		<b>0.05</b>	0.53	0.60
	ddG monomer (hard-rep)		-0.09	<b>0.52</b>	<b>0.73</b>
	ZEMu paper		-0.26	0.68	0.60
pdb-1CSE	flex ddG	6	0.53	1.83	0.67
	no backrub control		0.37	1.94	0.67
	ddG monomer (hard-rep)		0.46	1.94	0.67
	ZEMu paper		<b>0.87</b>	<b>1.06</b>	<b>1.00</b>
pdb-1DAN	flex ddG	118	0.65	<b>0.67</b>	<b>0.86</b>
	no backrub control		<b>0.69</b>	0.74	<b>0.86</b>
	ddG monomer (hard-rep)		0.61	0.71	0.83
	ZEMu paper		0.32	0.83	0.82
pdb-1DFJ	flex ddG	20	0.71	1.11	<b>0.70</b>
	no backrub control		<b>0.83</b>	<b>1.04</b>	0.65
	ddG monomer (hard-rep)		0.69	1.15	0.60
	ZEMu paper		0.55	1.26	0.60
pdb-1DQJ	flex ddG	34	<b>0.44</b>	<b>1.56</b>	<b>0.79</b>
	no backrub control		0.39	1.73	0.71
	ddG monomer (hard-rep)		0.37	1.66	<b>0.79</b>
	ZEMu paper		0.28	1.79	<b>0.79</b>
pdb-1DVF	flex ddG	38	<b>0.65</b>	1.21	0.61
	no backrub control		<b>0.65</b>	<b>1.12</b>	0.68
	ddG monomer (hard-rep)		0.61	1.27	<b>0.71</b>
	ZEMu paper		0.57	1.30	0.63
pdb-1E96	flex ddG	6	<b>0.80</b>	<b>0.76</b>	0.50
	no backrub control		0.51	0.83	0.50
	ddG monomer (hard-rep)		0.45	0.82	0.50
	ZEMu paper		0.50	0.78	<b>0.67</b>
pdb-1EAW	flex ddG	27	-0.02	0.57	0.85
	no backrub control		0.07	0.66	0.78
	ddG monomer (hard-rep)		<b>0.13</b>	<b>0.51</b>	0.89
	ZEMu paper		0.00	0.53	<b>0.93</b>
pdb-1EMV	flex ddG	51	<b>0.89</b>	<b>0.85</b>	<b>0.86</b>
	no backrub control		0.84	0.97	0.84
	ddG monomer (hard-rep)		0.84	0.98	0.84
	ZEMu paper		0.87	0.91	0.84
pdb-1F47	flex ddG	12	0.52	0.73	0.50
	no backrub control		0.58	0.67	<b>0.58</b>
	ddG monomer (hard-rep)		<b>0.60</b>	<b>0.66</b>	<b>0.58</b>
	ZEMu paper		0.51	0.75	0.50
	flex ddG	13	-0.15	0.84	0.56

Mutation Category	Prediction Method	N	R	MAE	FC
Complete dataset	flex ddG	1240	<b>0.63</b>	<b>0.93</b>	<b>0.76</b>
	flex ddG (REF energy)		<b>0.63</b>	<b>0.93</b>	<b>0.76</b>
Small-to-large mutation(s)	flex ddG	130	<b>0.64</b>	<b>0.87</b>	<b>0.75</b>
	flex ddG (REF energy)		0.57	0.92	0.72
Single mutation to alanine	flex ddG	748	<b>0.50</b>	<b>0.72</b>	0.75
	flex ddG (REF energy)		0.49	0.73	<b>0.76</b>
Multiple mutations	flex ddG	273	<b>0.62</b>	<b>1.51</b>	<b>0.77</b>
	flex ddG (REF energy)		0.59	1.57	0.75
Res. $\leq 1.5$ Ang.	flex ddG	52	0.46	0.85	0.75
	flex ddG (REF energy)		<b>0.65</b>	<b>0.74</b>	<b>0.77</b>
Res. $\geq 2.5$ Ang.	flex ddG	457	<b>0.50</b>	<b>0.74</b>	0.74
	flex ddG (REF energy)		0.48	0.75	<b>0.76</b>

Table 9: Performance comparison of the standard flex ddG protocol (using Rosetta’s Talaris energy function) with flex ddG run with the REF score function. "res  $\leq 1.5$  Ang." indicates data points for which the resolution of the input wild-type crystal structure is less than or equal to 1.5 Å. Backrub steps = 35000. R = Pearson’s R. MAE = Mean Absolute Error. FC = Fraction Correct. N = number of mutations in the dataset or subset.