

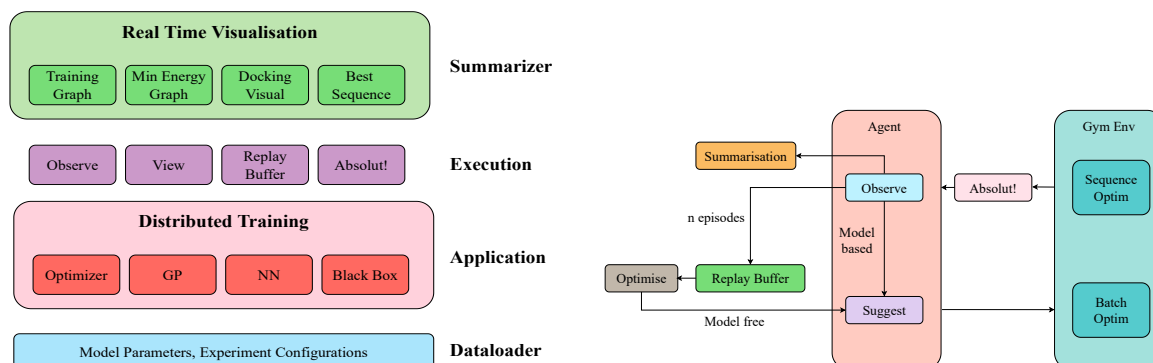
Cell Reports Methods, Volume 3

Supplemental information

**Toward real-world automated antibody design
with combinatorial Bayesian optimization**

Asif Khan, Alexander I. Cowen-Rivers, Antoine Grosnit, Derrick-Goh-Xin Deik, Philippe A. Robert, Victor Greiff, Eva Smorodina, Puneet Rawat, Rahmad Akbar, Kamil Dreczkowski, Rasul Tutunov, Dany Bou-Ammar, Jun Wang, Amos Storkey, and Haitham Bou-Ammar

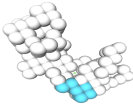

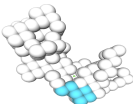
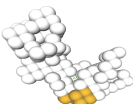

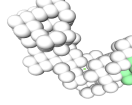
Supplementary Information



(a) Architecture of end-to-end framework for black-box optimisation. The architecture divides into four layers. The bottom layer consists of model parameters and experiment configurations which could be defined by the developers. The application layer, pre-written or written by developers, sets up the components for the execution layer. The detail of the execution layer is shown on the right side. The summarise layer collects data every iteration and produces real-time visualisation of the results.

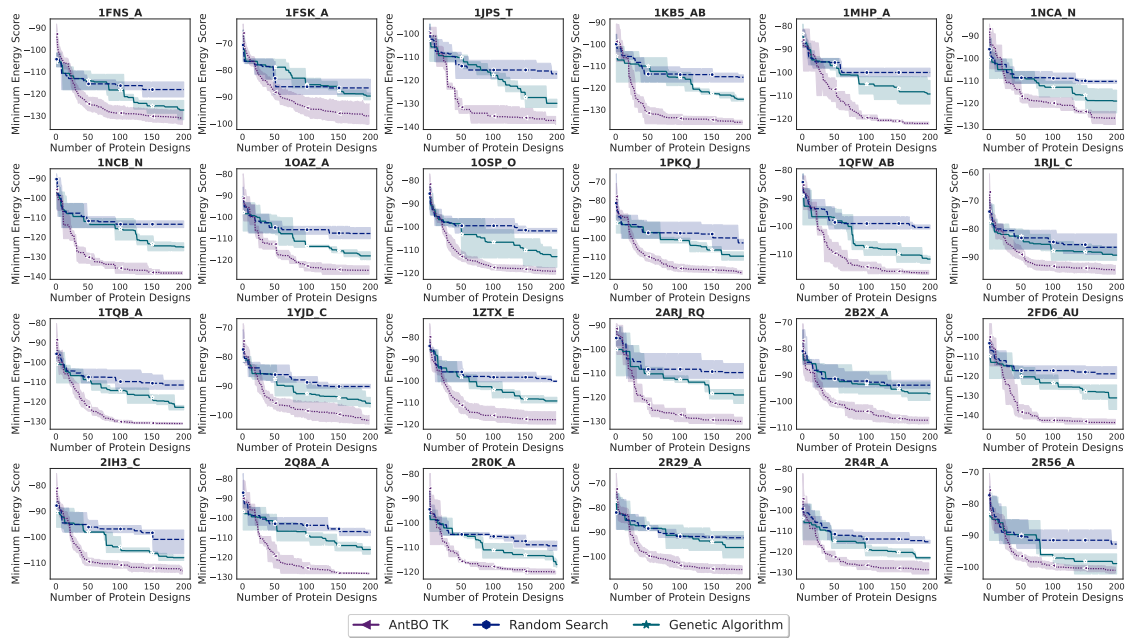
(b) Abstraction within the execution layer. The agent suggests the CDRH3 sequences and passes them into the gym environment. The gym environment evaluates the corresponding binding energies with Absolut!. The agent observes the results and calls the summarisation function to update real-time data. The results are stored in the replay buffer, which can be used to train deep reinforcement learning models. Within the observe function, the model-based agent also optimises the model. The agent then suggests the new CDRH3 sequences in the next iteration.

Figure S 1: Layout of AntBO software as introduced in Method Section 11.7 of a manuscript. On the left is the architecture of the framework. On the right is the illustration of the execution layer.

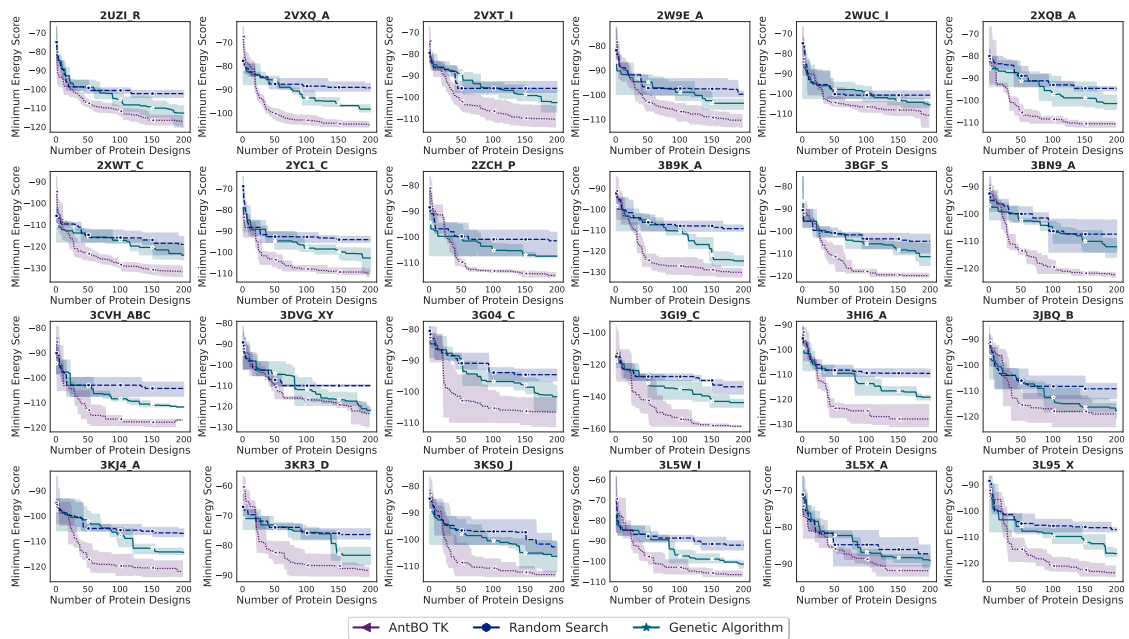
Step	Sequence	Energy	
0	HICAGFWHMPI	-88.31	
10	TDKTHPEVYTR	-67.36	
20	HYGMFLLPVGL	-97.67	
30	HGHMFLLHVIL	-90.44	
40	HFGMFELYVIL	-100.75	
50	HFVMFFLYVAL	-102.87	
60	HFVMPLLVLML	-98.53	
70	HFVMFYLVLMML*	-102.68	
80	HFTMFFLVLMML	-105.87	
90	HFTMCFLVLML	-102.31	
100	FFIMFFLQLIL	-106.61	
110	FFIMFFLVLCL	-109.5	
120	FFIMFPLVLI	-107.18	
130	PFIMFFLVLT	-104.53	
140	FFIMFLLVFL	-108.22	
150	FFIMFLLHLYL	-96.04	
160	CFIMFLLVLT	-107.29	
170	FFIFFLSVLWL	-102.17	
180	FFIFFLLFTIL	-107.2	
190	FFIFFVLFIL*	-110.42	

* +

Figure S 2: An example of a trajectory of sequences every ten steps generated by *AntBO*, annotated with their respective binding affinity, Related to Figure 2 and Section 3.4.1 of a manuscript. The structures of sequences are shown on the right. Each structure is denoted by a different colour, and from steps 40 to 197, the sequences share the same binding structure (in purple). Additionally, two sequences (70 and 190, marked with an asterisk) add an equally optimal binding structure (i.e., two binding modes), shown in green.

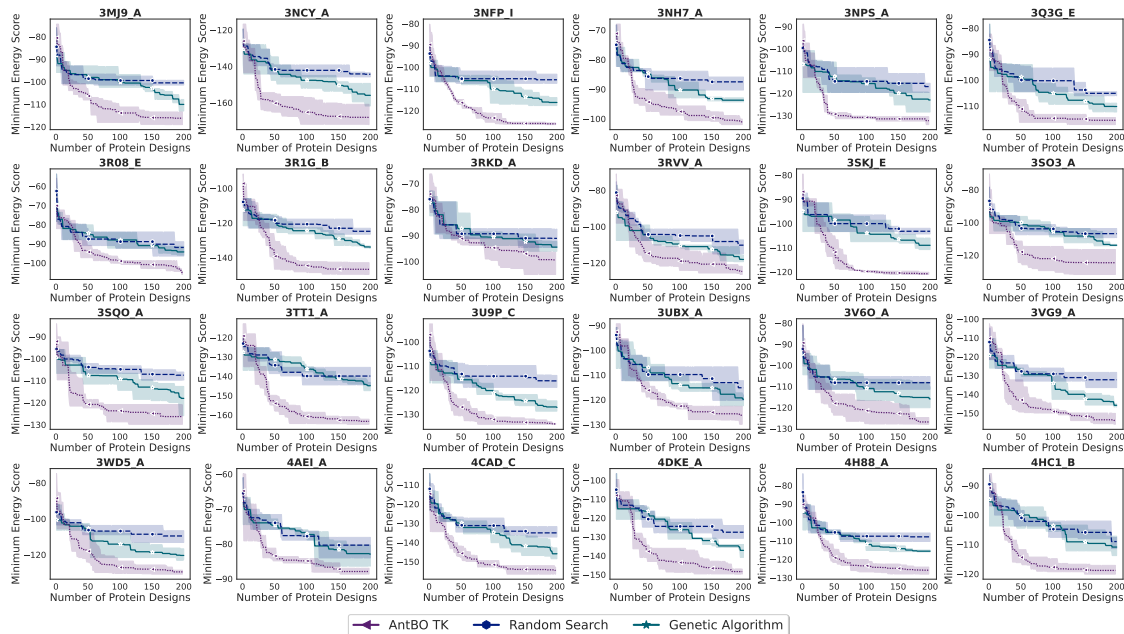


(a)

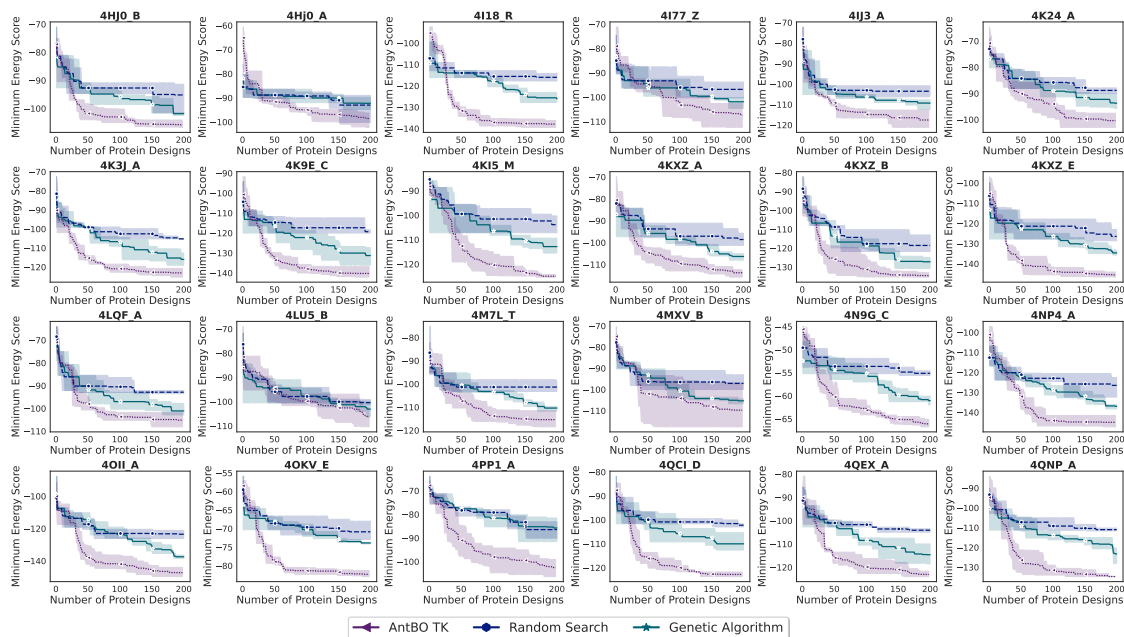


(b)

Figure S 3: Evaluation of AntBO on remaining antigens. Here, we report binding energy vs the number of protein designs comparing best performing AntBO TK with random search and genetic algorithm baseline. Related to Figure 2 of the manuscript.

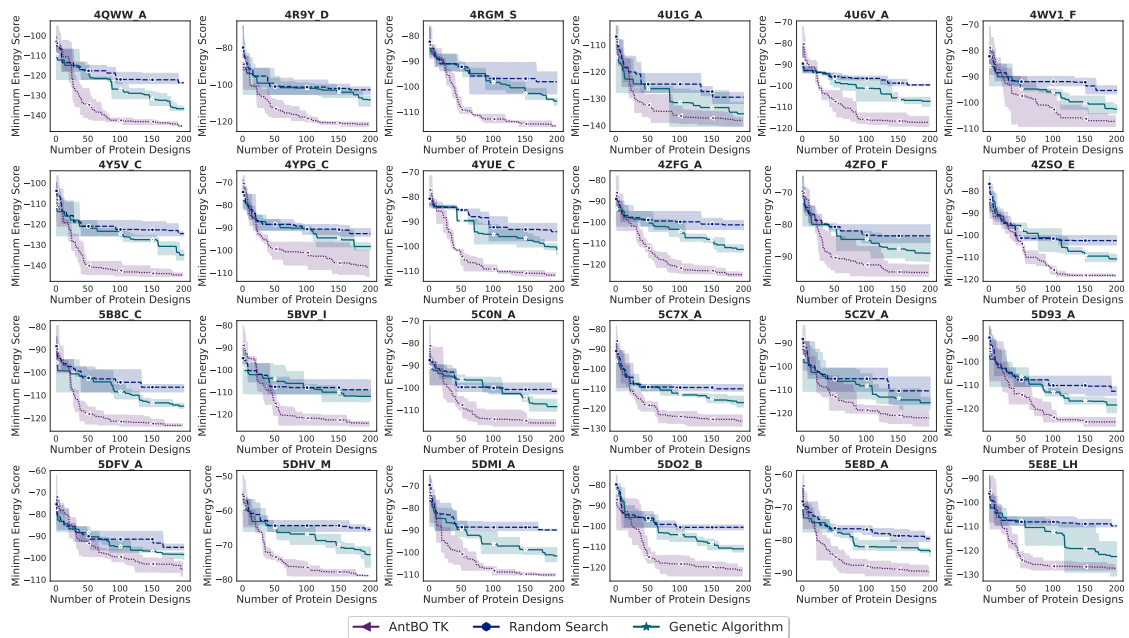


(a)

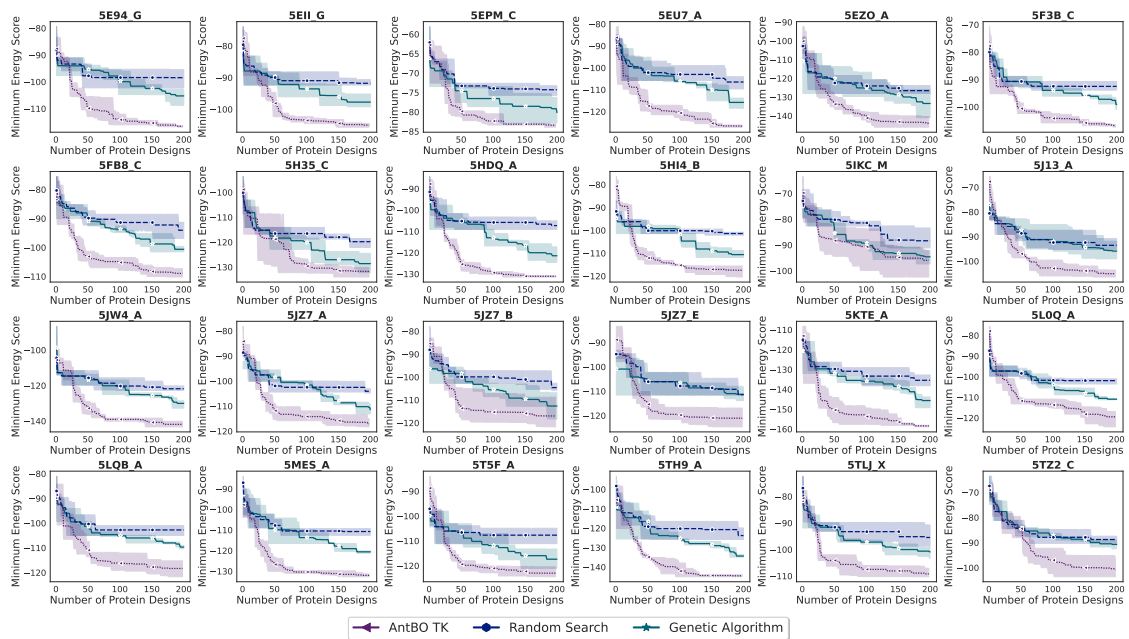


(b)

Figure S 4: Evaluation of AntBO on remaining antigens. Here, we report binding energy vs the number of protein designs comparing best performing AntBO TK with random search and genetic algorithm baseline. Related to Figure 2 of the manuscript.



(a)

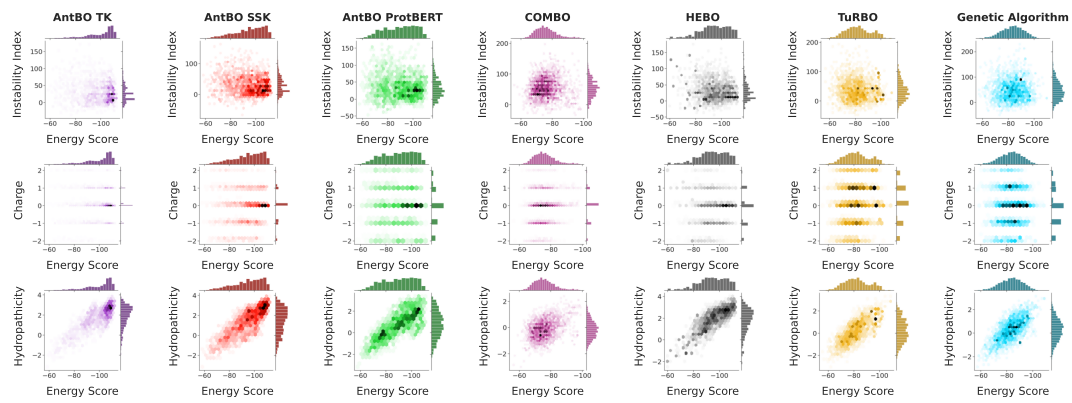


(b)

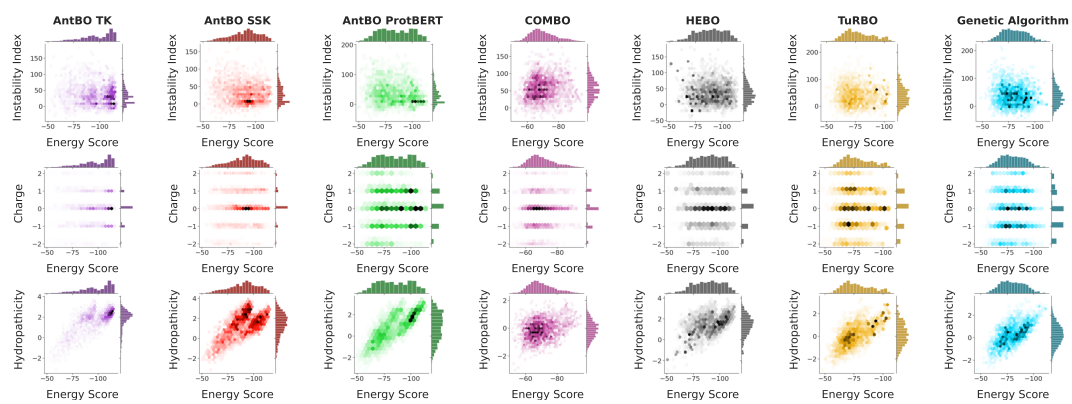
Figure S 5: Evaluation of AntBO on remaining antigens. Here, we report binding energy vs the number of protein designs comparing best performing AntBO TK with random search and genetic algorithm baseline. Related to Figure 2 of the manuscript.

Algorithm	Hyperparameter	Value
AntBO TK / AntBO NT / AntBO SSK / AntBO BERT	Acquisition function	Expected Improvement
	Nb. of initial points	20
	Normalise	True
	Kernel Type TK /	Transformed Overlap Kernel /
	Kernel Type NT /	Transformed Overlap Kernel /
	Kernel Type SSK /	Fast String Kernel /
	Kernel Type BERT	RBF Kernel with lengthscale on BERT features
	Noise Variance	1e-6
	Search Strategy	CDRH3 (trust-region) Local Search
	Use trust region TK / NT / SSK / BERT	Yes / No / Yes / Yes
COMBO	Trust region Length Min d_{\min}	1
	Trust region Length Max d_{\max}	30
	Batch size	1
	Nb. of initial points	20
	GP-parameters slice sampling steps	100 (init) / 10 (refine)
	Acquisition function	Expected Improvement
HEBO	Nb. of random samples for BFLS	20,000
	Nb. of initial points for BFLS	20
	Batch size	1
	Surrogate Model	Gaussian Process
	Acquisition Class	Evolution Optimiser
	Acquisition Optimiser	MACE
TURBO	Population Size	100
	Optimiser Nb. of Iterations	100
	Optimiser ES	NSGA-II
	Batch size	1
	Nb. of trust regions	1
	Trust region Length Min	2^{-7}
	Trust region Length Max	1.6
	Trust region Length Init	0.8
	τ_{succ}	3
	τ_{fail}	d
	Max Cholesky Size	2000
	GP fit - Optimiser	Adam
GP fit - Training Steps	50	
GP fit - Learning Rate	0.1	
Nb. of Thompson Samples	$\min(100d, 5000)$	
GA	Population size	40
	Nb. of iterations	5
	Nb. of parents	16
	Nb. of elite	6
	Crossover type	uniform
	Crossover probability	1.
	Elite ratio	0.15
	Mutation probability	$1/d$
RS	Nb. of iterations	200
	Sampling type	uniform
LamBO	Query batch size (b)	1
	Batch set size ($ \mathcal{X}_{\text{base}} $)	16
	Nb. of initial points ($ \mathcal{D}_0 $)	200
	Surrogate model	Exact GP (<code>single_task_exact_gp</code>)
	Acquisition function)	Expected Improvement
Encoder	<code>mlm_cnn</code>	

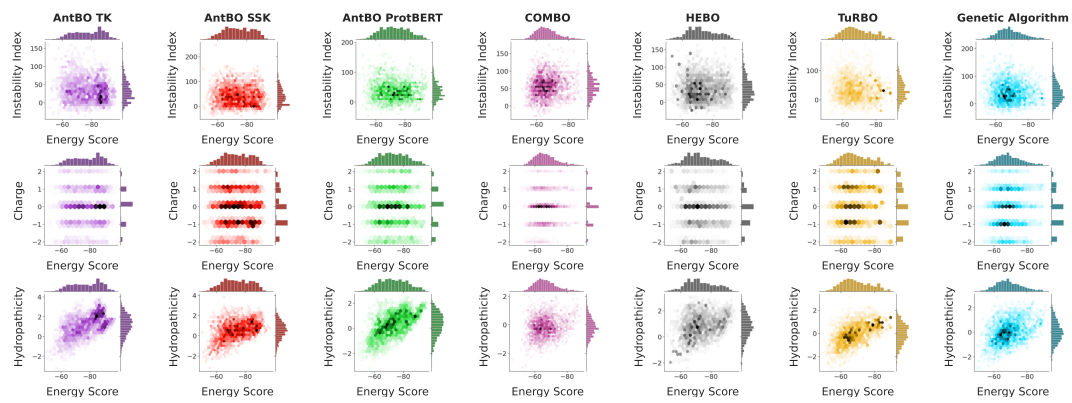
Table 1: Hyperparameter Configuration of different optimisation methods explained in the Method Section 11.8 of a manuscript.



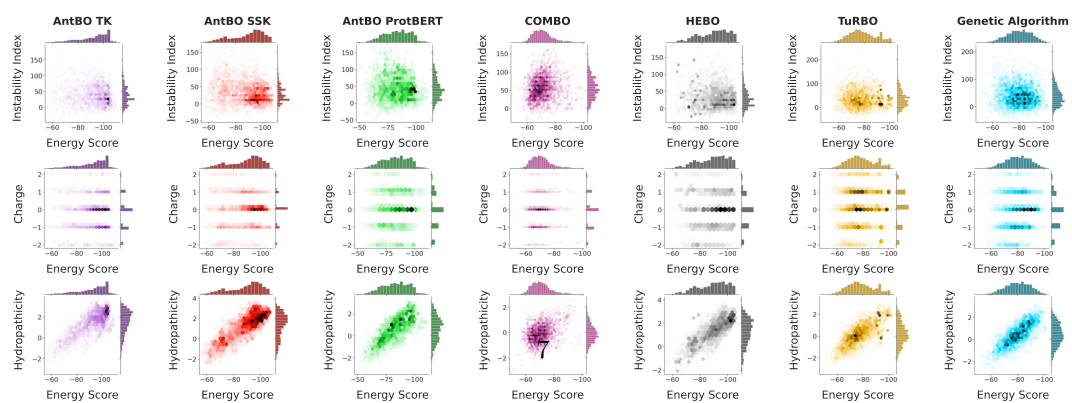
(a) 1ADQ (A)



(b) 1FBI (X)

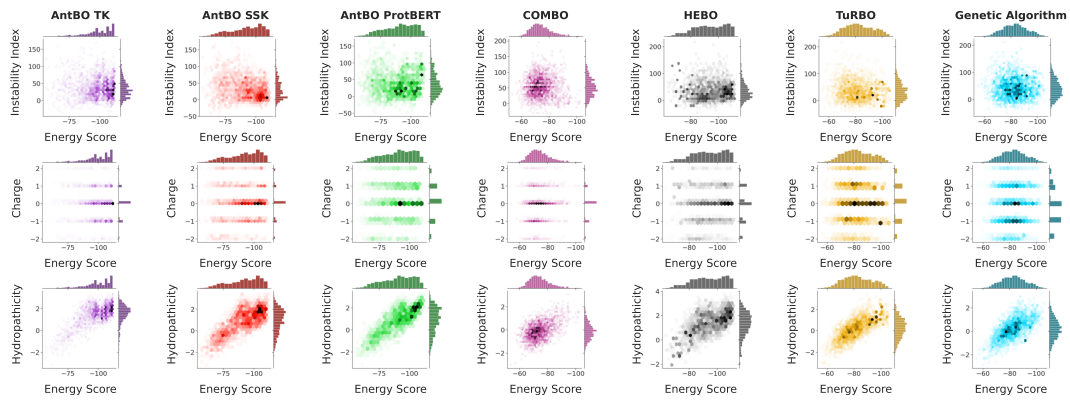


(c) 1H0D (C)

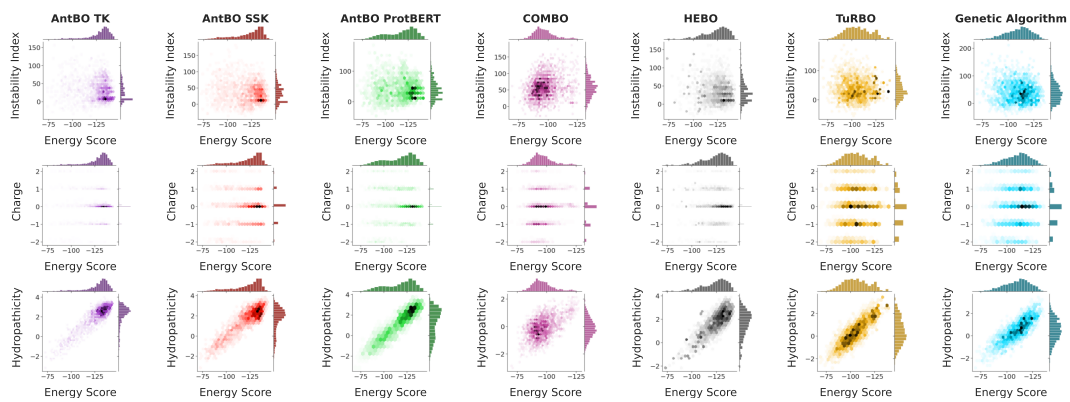


(d) 1NSN (S)

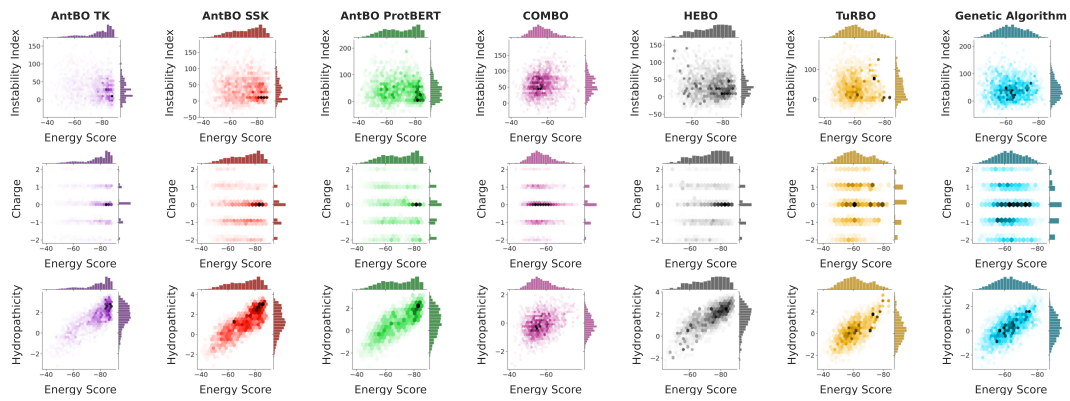
Figure S 6: We analyse the developability scores of 200 proteins designed by each method averaged across all 10 random seeds. Related to Figure 4 of the manuscript.



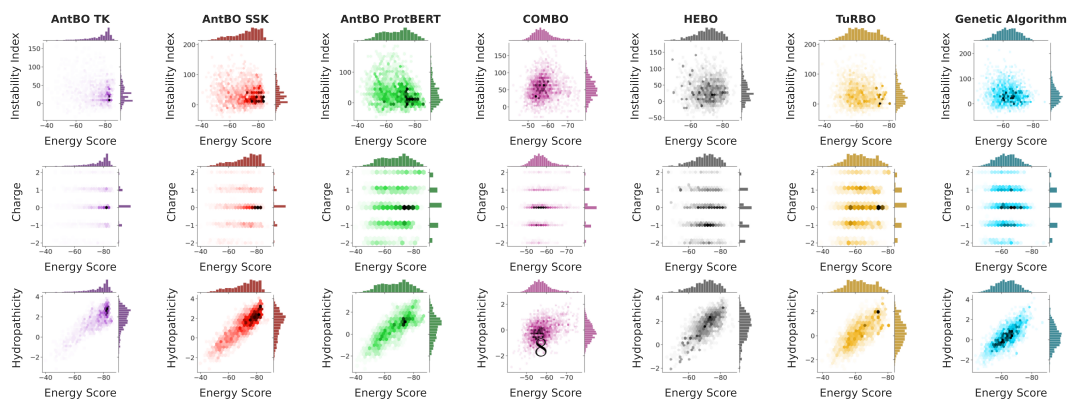
(a) IOB1 (C)



(b) 1S78 (B)

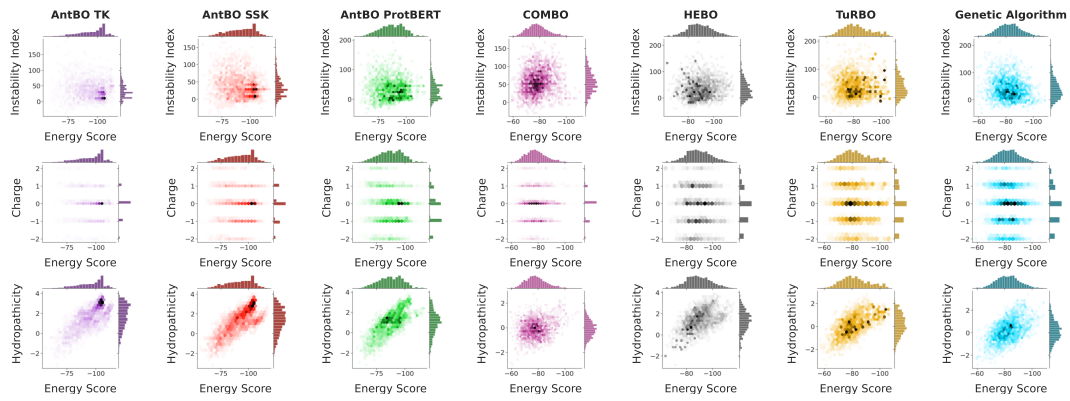


(c) 1WEJ (F)

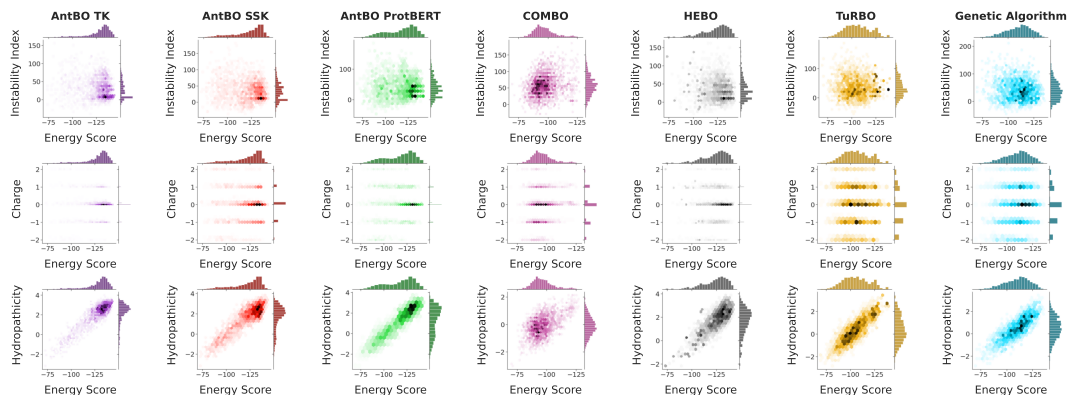


(d) 2JEL (P)

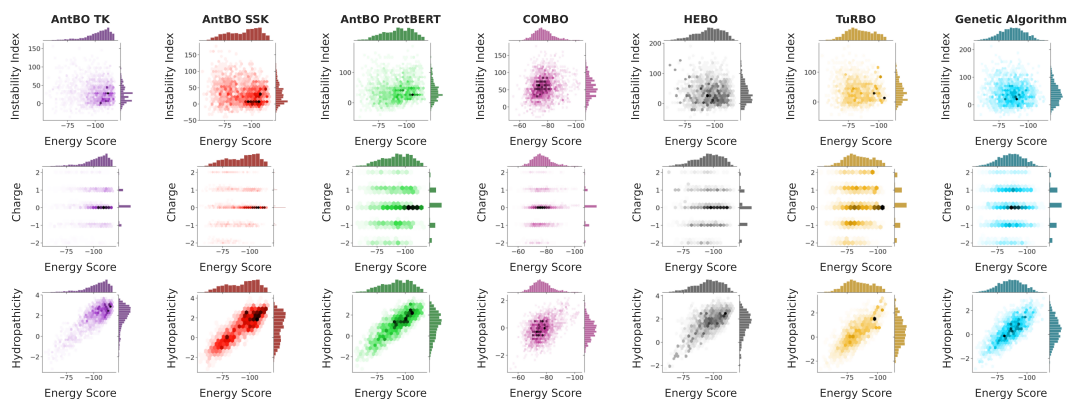
Figure S 7: We analyse the developability scores of 200 proteins designed by each method averaged across all 10 random seeds. Related to Figure 4 of the manuscript.



(a) 2YPV (A)



(b) 3RAJ (A)



(c) 3VRL (C)

Figure S 8: We analyse the developability scores of 200 proteins designed by each method averaged across all 10 random seeds. Related to Figure 4 of the manuscript.