

1 Thank-you for the thoughtful and constructive feedback on our manuscript. Below, citation numbers refer to the main  
2 text and new references are made in-line.

3 We have now run experiments showing that autofocus can still be beneficial in *higher dimensional design spaces*.  
4 For the superconductor design experiment, originally conducted on the 10 most informative features, we have run  
5 the same experiment on 60 features (leaving out the 21 least informative features due to colinearity issues). As a  
6 point of comparison, many Bayesian optimization methods are considered practically useful up to around ten or so  
7 dimensions (Wang, Zoghi, *et al.*, IJCAI 2013). Presently, 60 is already a stretch goal for many problems. We find that  
8 autofocus maintains its statistically significant gains over non-autofocusing for the trust region EDA, CbAS (Table A1)  
9 at dimensionality of 60. We have not yet examined the other methods, but this result alone already demonstrates that  
10 there is no inherent problem to using autofocus in higher dimensions. We will pursue more systematically investigating  
11 the effect of dimensionality, across all methods.

Table A1: As Table 1 in the main text, but with 60 dimensions. (\*) means  $p$ -value  $< 0.05$ , (\*\*) means  $p < 0.01$

	Median	Max	CbAS	$\rho$	RMS
Original	44.7	96.0		-0.12	26.4
Autofocused	52.0	103.4		0.12	21.6
Mean Diff.	7.2**	7.3**		0.23**	-4.7**

12 If one defines *extrapolation* as going into regions of the input space where the property (*e.g.*,  $\mathbb{E}[y | \mathbf{x}]$ ) takes on very  
13 different values from the training labels, then indeed, this does occur in our experiments. In particular, all autofocused  
14 methods except Random Search always produced *top candidates with ground-truth expectations greater than the*  
15 *maximum training label*. We have now computed the probabilities that the methods produce top candidates with  
16 ground-truth expectations that surpass the maximum training label (a “probability of improvement”). For example,  
17 non-autofocused CbAS yields a probability of improvement that is 28 times higher than that of a random control (*i.e.*,  
18 drawing the same number of samples from the training distribution) ( $p < 0.01$ , Mann-Whitney  $U$ -test). Autofocused  
19 CbAS increases this to 65 times ( $p < 0.01$ ). As for whether we are interpolating within the training data, we are not in  
20 the sense that the effective sample size (ESS) reported in the original submission shows that the search model deviated  
21 quite considerably from the training distribution (Fig. S3), as the ESS is the inverse of the exponentiated Renyi-2  
22 divergence. Similarly, in additional experiments, we now find that the KL divergence between the training and test  
23 distributions goes up to 150 for DbAS and up to 10 for CbAS. All that being said, the goal of AF is a subtle one. We  
24 do not claim that AF *necessarily* makes a design method move further beyond the training data; instead, AF helps  
25 us make the best use of the oracle for wherever a design method does move, regardless of whether that movement  
26 is interpolation or extrapolation (also see “Oracle bias-variance trade-off” paragraph on page 4). From that point of  
27 view, using a method with a trust region complements, rather than contradicts, the goal of AF, although a more formal  
28 examination of how trust-region methods like CbAS can still enable extrapolation would be useful.

29 With respect to “reporting the expected probability of belonging to  $S$ , when say  $y_\tau$  is the 80th percentile”, if we  
30 understand correctly, this is tantamount to asking about extrapolation, and hence addressed in the previous paragraph.  
31 Note too that for maximization design problems, the desired  $S$  will be the set of values  $y$  such that  $y \geq y_{\max}$  (where  
32  $y_{\max} \equiv \max_{\mathbf{x}} \mathbb{E}_{p(y|\mathbf{x})}[y]$ ).

33 With respect to *discrete design spaces*, in preliminary experiments on the protein design problem in (9), autofocused  
34 CbAS compared to CbAS yielded candidates with a maximum ground truth of 3.39 compared to 3.37 ( $p = 0.03$ ) and a  
35 Spearman’s  $\rho$  of 0.83 instead of 0.70 ( $p = 0.01$ ). This topic requires more comprehensive experimentation which we  
36 continue to do.

37 The practical importance of the *fixed-data setting* has been demonstrated by recent protein engineering and materials  
38 design work, in which a single labeled dataset is collected, and a regression model is trained and used to guide the  
39 design of novel proteins without additional data (5, 6). Moreover, a prominent protein engineer has told us the fixed-data  
40 setting is what they are currently working on (Frances Arnold, personal communication). Finally, no matter how many  
41 rounds of data are acquired, at some point, the data are fixed and our current problem formulation applies.

42 With respect to the comment about non-negligible  $P_0(y \in S)$ , indeed, if this is too small, then the variance explodes. In  
43 effect, this says that the starting (training) data must have at least a hint of the property we care about, just as is required  
44 for successful directed evolution for protein design (3). We will add further clarification in the manuscript.

45 We will incorporate feedback on clarity of Fig 1, and also on our related work section, in particular linking to  
46 out-of-distribution generalization, small-area estimation, model misspecification, model inversion networks, and the  
47 Angermueller paper *etc.* As for formulating autofocus as a game, this was simply done to give a precise name to  
48 the problem formulation that emerged from our derivation. And indeed, Nash equilibria are hard, but as you note,  
49 empirically, we still find benefit from our approach. Finally, we apologize if some points were not addressed, but in  
50 light of time and space constraints we had to prioritize. We will of course address all points raised.