

1 We thank the reviewers (R1, R2, R3) for the helpful comments, corrections and suggestions.

2 The main concern seems to be the **robustness** of our results to various deviations from the idealized scenario considered
3 by our theory (R2, R3). First, while the **hyperparameter** space is too large to explore systematically, simulations
4 suggest the qualitative phenomenology presented in the paper is robust to various model details. The effects of varying
5 the fraction of active/informative neurons are shown in Fig. 2B/3C (we will improve 2B to include all decoders
6 (R2)). We will also document additional parameter variations in the supplementary material (SM). Second, R2 and R3
7 expressed concern about the precision required for the **optimal modulation weights** in the encoding. Here we show
8 numerically that the results hold qualitatively even when the modulation deviates significantly from the absolute optimal
9 decoding weights ($w = dec_{ML} + \epsilon$ where ϵ is independent gaussian noise; example in panel A) although the overall
10 performance degrades (B). Hence our idea could still apply to a more realistic suboptimal encoding model.

11 R2 also identified several potential mismatches between our idealized model and real data. First, the modulator could
12 be **multidimensional**. We chose to focus on the unidimensional case because: 1) it is the simplest, 2) mathematically,
13 having a linear combination of gaussian, task-specific modulators does not qualitatively change the problem, though
14 it makes the model harder to parametrize (additionally, if some modulator targeting is not task specific, that would
15 reduce the SNR of all neurons and correspondingly affect all decoders), and 3) while in the Rabinowitz paper there
16 were several modulators, most of variance was accounted for by one (for each hemisphere). Second, the presence of
17 **additive noise**: experimental reports are conflicting (e.g. Goris et al.2014 argue that additive noise is inconsistent with
18 their data); moreover, to date, there is no evidence that this component is functionally targeted. Simulations show that
19 task-invariant additive noise decreases the performance of all decoders, but does not qualitatively change our results
20 (panel C, to be included in SM). These issues will be further detailed in the Discussion.

21 We chose to focus on classification rather than **estimation** because the experiments showing task-specific modulation
22 use binary discrimination. As R2 rightfully points out, it is important to expand the theory to other tasks. In principle,
23 since estimation also entails learning to appropriately weight informative neurons while ignoring uninformative ones,
24 modulator-labeling should be helpful there too, though the details of the best encoder and decoder will likely change.
25 Including an informative **prior** (R2) should not qualitative change the discussion (it only shifts the threshold), unless the
26 prior directly affects the pattern of modulation. To our knowledge, there is no experimental data supporting this idea.

27 The model makes several **experimental predictions** (R2), which we are in the process of testing using V1 monkey
28 recordings in a task that dynamically shifts the task-relevant sub-population. The main assumptions of the theory
29 to check: 1) the subpopulation of task-informative neurons is small and hence hard to identify, 2) low-dimensional,
30 shared noise, changing faster than the trial duration, that preferentially targets informative neurons. We further predict
31 that our modulator-guided heuristic decoder should outperform simpler strategies (sign-only or rate-guided). Perhaps
32 counterintuitively, **attention** reduces the variance of the modulator (Rabinowitz et al). Since our theory posits an
33 optimal level of modulation relative to the stimulus-induced variance, we would suggest that attention shifts the level
34 of modulation towards this (empirically estimatable) optimum. The direction of the shift may provide hints about the
35 mechanics of noise generation, something we know little about at the moment (Huang et al. 2019).

36 There is a small misunderstanding regarding the usage of **MSE** (R1), which we use only for the *modulation weights*
37 estimator, but not for assessing task performance (measured as % correct). Similarly, the procedure used for the **learning**
38 of the MG decoding weights seems unclear (R2). It includes two processes: 1) learning the signs, which happens using
39 end-of-trial feedback, but only needs few examples (see Fig.2A), and 2) learning the absolute optimal decoding weights,
40 which happens within the trial by estimating modulation weights, with variance scaling as $1/T$; the exact learning
41 rate is determined by the strength and time constant of the modulator (see Eq.14). Quantitatively confirming if the
42 modulation is enough to explain the animal's **speed of learning** requires further data analysis (ongoing).

43 Minor: the explicit reference of the target behavioral 90% level is somewhat misleading and will be removed (R2). We
44 will add the additional references (R2), clarify the equations (R2, R3) and improve Fig.3B visually and for clarification
45 (vertical line=optimum from 3A, middle panel should show precision instead of variance, to reduce confusion) (R1, R2).
46 All the above points will be incorporated in the camera-ready version of the paper, should our submission be accepted.

