

1 **The Hard Limits of Decoding Mental States: The Decodability of fMRI**

2
3 R. Jabakhanji^{1,8†}
4 A.D. Vigotsky^{2†}
5 J. Bielefeld^{1,8}
6 L. Huang^{1,8}
7 M.N. Baliki^{3,4,8}
8 G.D. Iannetti^{5,6}
9 A.V. Apkarian^{1,3,7,8*}

10
11 ¹ Department of Physiology, Feinberg School of Medicine, Northwestern University, Chicago, USA.

12 ² Departments of Biomedical Engineering and Statistics, Northwestern University, Evanston, USA.

13 ³ Department of Physical Medicine and Rehabilitation, Feinberg School of Medicine, Northwestern University,
14 Chicago, USA.

15 ⁴ Shirley Ryan AbilityLab, Chicago, USA.

16 ⁵ Division of Biosciences, University College London, London, UK.

17 ⁶ Neuroscience and Behaviour Laboratory, Italian Institute of Technology, Rome, Italy.

18 ⁷ Department of Anesthesiology, Feinberg School of Medicine, Northwestern University, Chicago, USA.

19 ⁸ Center for Translational Pain Research, Feinberg School of Medicine, Northwestern University, Chicago, USA.

20 * Correspondence to: a-apkarian@northwestern.edu.

21
22 † These authors contributed equally to this work.

23
24
25 **SUMMARY (150-word max)**

26 High-profile studies claim to assess mental states across individuals using multi-voxel decoders of brain
27 activity. The fixed, fine-grained, multi-voxel patterns in these “optimized” decoders are purportedly necessary
28 for discriminating between, and accurately identifying, mental states. Here, we present compelling evidence
29 that the efficacy of these decoders is overstated. Across a variety of tasks, decoder patterns were not necessary.
30 Not only were “optimized decoders” spatially imprecise and 90% redundant, but they also performed similarly
31 to simpler decoders, built from average brain activity. We distinguish decoder performance when used for
32 discriminating between, in contrast to identifying, mental states, and show even when discrimination
33 performance is strong, identification can be poor. Using similarity rules, we derived novel and intuitive
34 discriminability metrics that capture 95% and 68% of discrimination performance within- and across-subjects,
35 respectively. These findings demonstrate that current across-subject decoders remain inadequate for real-life
36 decision making.

38 INTRODUCTION

39 A whole body of neuroimaging literature—largely published in highly influential journals—either
40 explicitly claim, or strongly imply, that thinking is no longer private. By “optimizing” functional magnetic
41 resonance imaging (fMRI) brain scan results, these studies profess to universally decode mental states: feelings,
42 thoughts, decisions, intentions, and behaviors¹⁻³. Thus, neuroscience seems to have broken the code of mental states,
43 in turn proclaiming the ability to “read the brain” of every human being. Here, we systematically examine the
44 validity of such claims.

45 Decodability—how discernable a mental state is, given a brain activity pattern—is predicated both on the
46 brain activity properties of the task being discerned as well as the goal of the decoding. Intuitively, decodability is
47 analogous to discerning a breed of dog; breeds that look more similar will be harder to distinguish. The literature
48 claims “optimized” decoders can (1) discriminate between mental states, (2) identify mental states, and (3) capture
49 additional state-related measures (stimulus or perception intensities). These goals can be more tangibly elucidated
50 through the dog breed metaphor: Consider a pug (a *decodee*) and a French Bulldog (a *comparator*)—two breeds that
51 may look alike. If one is familiar with the unique physical features of a pug—small stature, short snout, wrinkled
52 face, folded ears, curled tail, etc.—then such features can serve as the *decoder* for a pug. This decoder can then be
53 used to perform the three decoding tasks. Specifically, discrimination is akin to deciding which dog is a pug when
54 the pug and French Bulldog are next to one another. Conversely, identification is akin to saying whether a single dog
55 is a pug when there are no other dogs around. On the other hand, capturing a continuous measure, such as perceived
56 intensity of a state, is much like trying to judge a dog’s age. While discrimination and capturing continuous
57 measures have been discussed and illustrated for various mental states, less attention has been given to identify a
58 certain mental state from a given pattern of brain activity.

59 The pattern of mental state decoders arises from weights assigned to its constituent voxels; for this reason,
60 we call them *fixed-weight decoders*. Voxel weights are derived in three stages. First, general linear models (GLM)
61 generate a brain *activity map* (correlation between the activity in each voxel and the task). Second, GLM is used to
62 contrast the activity maps from a task or state of interest (a *decodee*; e.g., pain) to one of no interest (a *comparator*;
63 e.g., touch), and its results are thresholded (a *contrast map*). Finally, “machine learning” models are used to tune the
64 weights in the contrast map to optimize its performance^{4,5}; the result is a relatively sparse, fixed-weight decoder
65 with a fine-grained pattern (an “*optimized*” *decoder*). It is tacitly assumed that each stage improves performance of
66 the decoder by uncovering better distributed patterns of neural ensembles related to the mental state, and as a result,
67 detailed spatial patterns confer predictive value, as explicitly stated, “the pattern of activation, rather than the overall
68 level of activation of a region, is the critical agent of discrimination”⁵. This concept is now expounded for diverse
69 topics across many labs^{3,5-11}.

70 The concept that across-subject “optimized” decoders are able to capture mental states across different
71 individuals violates basic neuroscientific principles. The technical and biological requirements of such decoders are
72 quixotic, as they imply the existence of a fixed, exclusive, universal brain activation pattern for each and every
73 mental state—a one-to-one correspondence between subjectivity and objective brain patterns. Such invariant brain-
74 to-mind models imply a common neuronal firing pattern across billions of neurons, which is unique for every mental
75 state and shared across all humans. This invariance is purported in spite of large inter-subject variability in gross
76 brain anatomy, as well as of differences in genetics, lifestyles, lifetime experiences, and associated memory traces
77^{12,13}; all of which would carve the individualized brain activity of subjectivity (for a discussion on the topic from the
78 viewpoint of fMRI analysis, see¹⁴). If a trivial, fixed relationship exists between subjectivity and brain activity, such
79 “optimized” decoders also raise strong ethical and legal concerns regarding their ability to invade mental privacy¹⁵,
80 and also would be incongruent with commonly accepted philosophical constructs of subjectivity¹⁶.

81 Our principal aim was to evaluate the performance and necessity of “optimized” decoders relative to more
82 parsimonious approaches (e.g., GLM maps). After rigorously evaluating the performance of “optimized” decoders,
83 we sought to understand fixed-weight decoders from a more general perspective: What determines and constrains
84 decodability?
85

86 RESULTS

87 *Overview*

88 Our investigation began with three published pain decoders. Both qualitatively and quantitatively, these decoders
89 were markedly different from one another (**Fig 1**). Despite these differences, on average, their ability to discriminate
90 pain from non-pain states, across four published studies ($N=113$)^{4,5,8,17}, was nearly identical. To understand how
91 disparate decoders could perform similarly, we parametrically perturbed each of the decoders and tested its
92 performance. Perturbations consisted of 1) searching for brain locations privileged for decoding pain; 2) randomly
93 using subsets of voxels from select regions; 3) using subsets of voxels based on their weights; and 4) spatially
94 smoothing (and thereby modifying voxel weights). The analysis demonstrated that the tested decoders were ~90%
95 redundant in space and, remarkably, that their weights were superfluous for successful discrimination (**Fig 2**).
96 Similar results were obtained for stimulus-perception mapping (**Fig 3**). Overall, we observed that sparse, location-
97 only based decoders were sufficient for discriminating pain.

98 To further generalize this finding, we examined decoding properties for cognitive domains other than pain,
99 where dedicated brain tissue is better established; namely, a reading task and a listening task (two publicly available
100 datasets, $n=14$ and $n=213$ subjects, respectively)^{18,19}. We compared decoding performance between GLM and
101 “optimized” decoders, before and after constraining to location-only. Our results closely resembled those for
102 decoding pain (**Fig 4**).

103 The brain imaging literature commonly accepts that if a decoder can adequately discriminate between a
104 decodee and a comparator, then it is also useful for identifying the mental state associated with the decodee. We
105 tested this concept for both pain and listening tasks. Despite discrimination being possible and robust to
106 perturbations, all decoders performed poorly and relatively similarly when trying to *identify* the decodee mental state
107 (**Fig 5**).

108 The results of our perturbation analyses led us to explore the limits of decoding. If perturbed and simplified
109 decoders can perform similarly to the original “optimized” decoders, can we further simplify decoders and also
110 quantitate decodability? To address the former question, we built pain decoders using GLM maps for noxious
111 stimuli. These GLM decoders performed similarly to “optimized” decoders, with within-study performance being
112 slightly superior to across-study performance (**Fig 6a,b**). We extended these findings to quantify within- and across-
113 subject decoding using four different tasks, repeated up to 12 times per subject in 14 subjects¹⁹. This study design
114 provides the opportunity to calculate discriminability as a function of similarity measures from the decoder,
115 decodee, and comparator, for both within- and across-subject decoding. Although performance was not consistently
116 better for within-subject discrimination, variation in performance could be largely explained by within-task
117 homogeneity and between-task heterogeneity, allowing us to propose decoding rules (**Fig 6c,d**), which worked
118 better for explaining within- compared to between-subject discriminability. These results present convergent
119 evidence that 1) discrimination decoding is limited by GLM results, where sparse location-only maps contain
120 sufficient information; 2) identification is harder than discrimination; 3) similarity measures almost fully account for
121 the variance of within-subject discrimination decodability, which degrades in across-subject discrimination.

122 *Exploring Established Decoders*

124 We started by assessing “optimized” decoders for pain using the Neurologic Pain Signature (NPS, constructed using
125 LASSO-PCR⁵) and Pain-Preferring Voxels (pPV, constructed using SVM⁴). In addition, we used a meta-contrast
126 map as an alternative decoder, Pain-Neurosynth (pNsy), which is the meta-analytic association test result for the
127 term “pain”. This contrast map is based on 516 studies containing the word “pain” in the abstract, and contrasting
128 them with the remaining 13,855 neuroimaging studies (using the public tool Neurosynth²⁰; see Supplementary
129 Methods). We first compared the spatial and weight properties of these three decoders. Although all decoders cover
130 approximately the same brain regions (**Fig 1A**), the distributions of their voxel weights are distinct (**Fig S1**), the
131 numbers of their constituent voxels are vastly different, their pairwise correlations are weak, and their spatial
132 overlaps of voxels are relatively small (**Fig 1, B-D**).

133

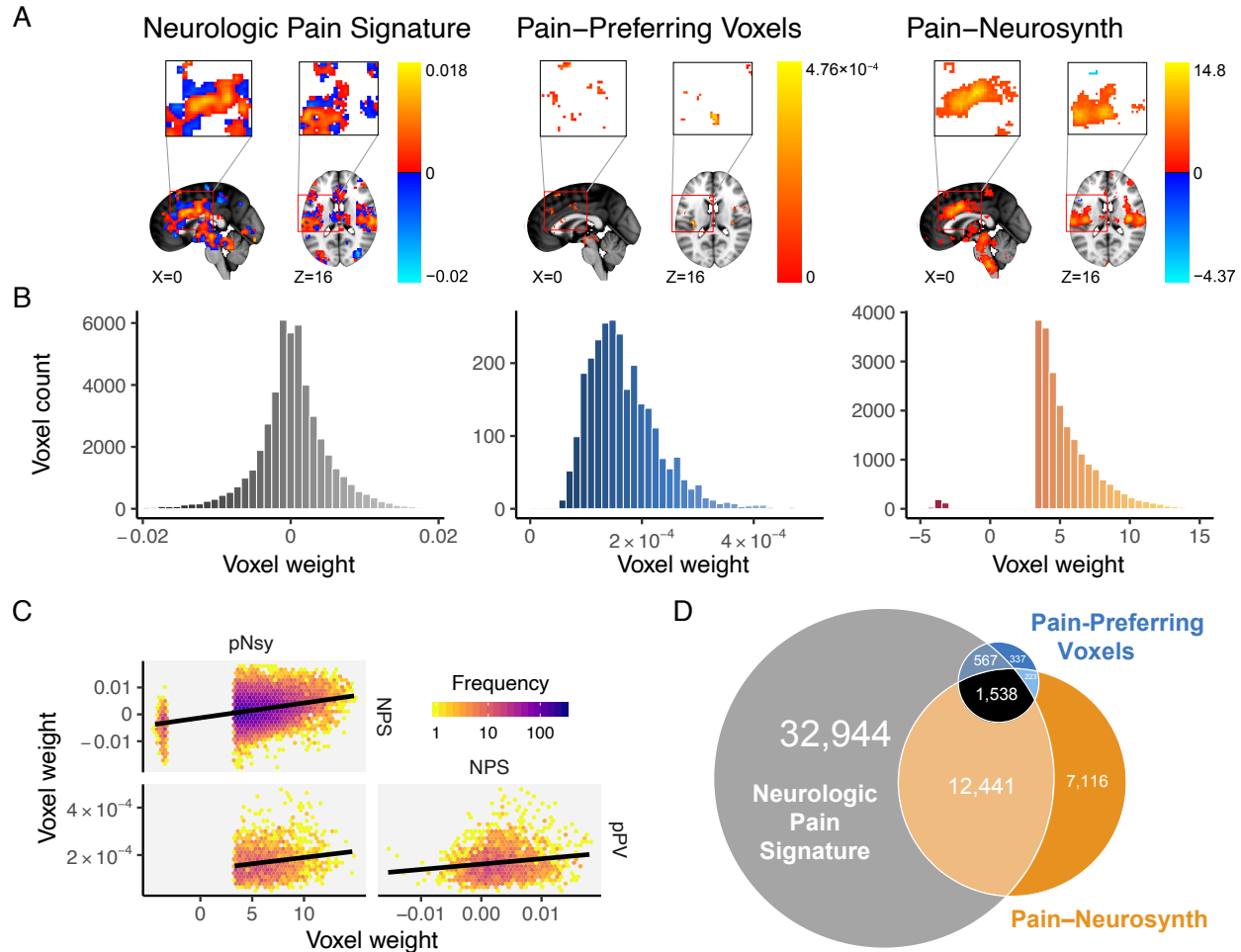


Figure 1. Spatial properties for three decoders, which are supposed to distinguish pain from other mental states, are distinct from each other. (A) Location and voxel-wise weight patterns of the three pain decoders (respectively abbreviated NPS, pPV, and pNsy). (B) Weight distributions of all three decoders are distinct. NPS weight values are distributed around zero; pPV has no negative weights; pNsy has only a few negative weights. (C) Pairwise correlations between weights of the three decoders. Lines depict total least squares regression fits. All three correlations are weak ($r_{\text{NPS-pPV}} = 0.16$; $r_{\text{pNsy-NPS}} = 0.30$; $r_{\text{pNsy-pPV}} = 0.18$). (D) Euler diagram depicts relative size of each of, and spatial overlap between, the three decoders.

Discrimination Performance for Pain is Similar Between Diverse Decoders

To enable decoding, we assessed the similarity between the decoder and decodee or comparator using the normalized dot product (NDP; +1 indicates total/maximal pattern similarity, 0 indicates orthogonal patterns, -1 indicates anti-similarity). Discrimination was assessed using the area under the receiver operating characteristic curve (AUC) from the two distributions of NDPs. AUC is an indicator of discriminability since it can be interpreted as the probability of a randomly sampled decodee NDP being greater than a randomly sampled comparator NDP, implying a direct comparison. Conversely, identifiability was assessed using distributional overlap, with greater overlaps indicating poorer identifiability. Points contained within the area of overlap are equally likely to be in the decodee and comparator distributions, and thus, are not identifiable. Together, these metrics served as the basis for decoding performance throughout this study.

We examined the performance of the three pain decoders for discriminating pain states from control tasks, and for capturing stimulus/perception properties, in 4 published studies from 3 labs ($N=113$ subjects). Despite marked spatial and weight distribution differences, average discrimination performance (meta-analysis, AUCs pooled across datasets and various comparators) were approximately equivalent ($\text{AUC} \approx 0.73$ for all three; **Fig 2A**).

157 This equivalence is remarkable and informative: it implies that very different models may nonetheless yield similar
158 average performance, suggesting that their detailed properties do not constrain decodability. Notwithstanding similar
159 average performance, the decoders performed differently for particular datasets, indicating that decoding
160 performance has a specificity component which can likely be explained by brain region-specific dependences.

161

162 ***Pain Decoders Are Robust to Spatial Perturbations***

163 *Search for Brain Locations Privileged for Decoding Pain*

164 To test whether there are privileged locations for decoding pain, we created clusters of voxels using the common
165 space across the three pain decoders and tested discrimination performance for all three decoders across the four
166 tasks (**Fig S2**). For any given study, multiple clusters from multiple decoders performed equally well and matched
167 the performance of the full decoder. This result suggests that no single cluster was consistently more specific for
168 decoding pain than other clusters.

169

170 *Spatial Smoothing and Voxel Weights*

171 To investigate whether discrimination performance relies on the *fixed-weight* nature of the voxel patterns, we
172 measured performance when these patterns were degraded (1) by spatial smoothing and (2) by discarding their
173 weights. Remarkably, decoding performance was minimally affected by either procedure (**Fig 2B, Figs. S3-S6**).
174 This result clearly demonstrates that the fine-grained pattern of weights in “optimal” decoders added no value to
175 performance (with a few exceptions, **Fig S3**); rather, voxel locations alone were sufficient for discrimination.

176

177 *Number of Voxels*

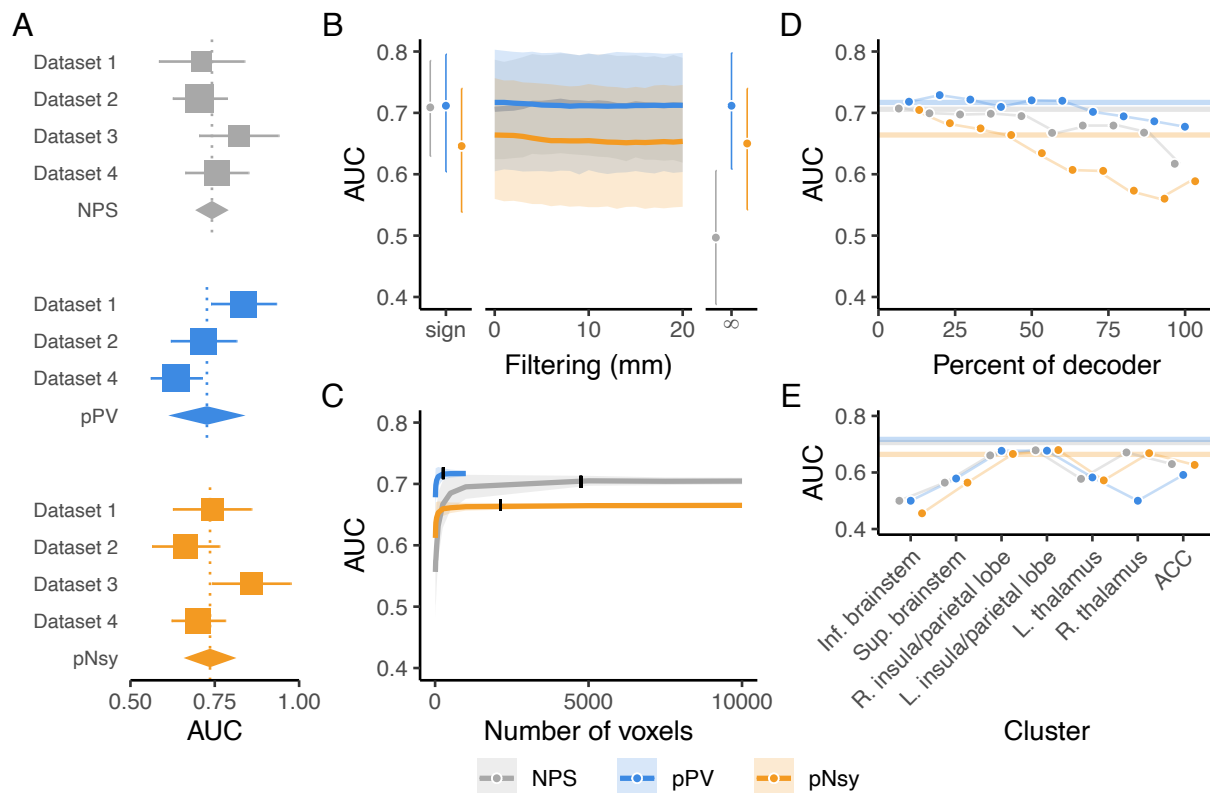
178 Given the three decoders’ structural differences and yet similar performances, it can be virtually ruled out that all
179 voxels composing these decoders are necessary to achieve satisfactory discrimination. To characterize the minimum
180 number of voxels necessary to decode the pain state, we created sets of new decoders by randomly selecting subsets
181 of voxels from each decoder. Surprisingly, we attained the original decoding performance when only using a
182 random 10% of the total number of each decoder’s constituent voxels (**Fig 2C**). We replicated this finding on all
183 datasets and for all three decoders, both in their original form and when only using their signed voxel locations
184 (**Figs. S4-S6**). We further explored the relationship between voxel weights and performance by first binning voxels
185 by their absolute weights and then constructing a set of decoders using the voxels in each bin (see **Fig S7**). Again,
186 we observed only a minimal degradation of performance with decreasing voxel weights for all decoder-dataset
187 combinations (with some degradation seen mainly for unfiltered NPS at highest binning, **Figs. S8-S9**). Together,
188 these results demonstrate that the information within the set of voxels present in the decoder, decodee, and relative
189 to the comparators was highly redundant and essentially independent of the decoders’ weights.

190

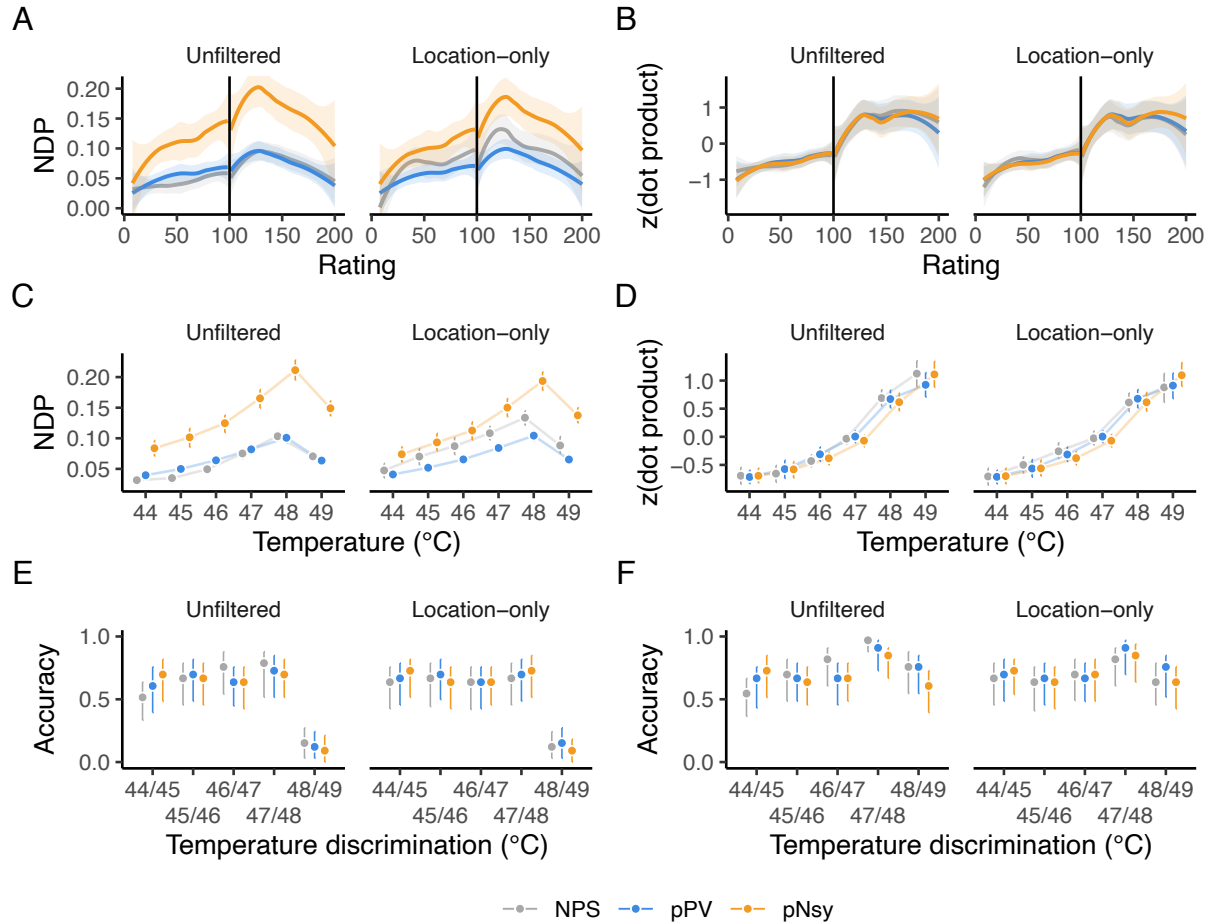
191 *Fixed, Multi-voxel Patterns Confer No Added Value to Stimulus/Perception Intensity Decoding*

192 Wager, et al. ⁵ used an “optimized” decoder, NPS, not only to discern the dichotomous presence of pain, but also to
193 claim that NPS can capture stimulus intensity and perceptual ratings from brain activity. To this end, we tested the
194 ability of the three pain decoders to capture stimulus and perception properties. We used data from a study where
195 nonpainful and painful stimulus, perceptual responses, and their associated brain activity were available ⁵. All three
196 decoders (NPS, pPV, and pNsy), whether unfiltered or infinitely filtered (location-only), performed similarly for
197 capturing perceived pain ratings (**Fig 3A,B**), for reflecting the intensity of the thermal stimulus (**Fig 3C,D**), and for
198 discriminating between pairs of painful stimuli (**Fig 3E,F**). The similar performance between unfiltered and
199 location-only decoders demonstrates that response trends arise as a consequence of the intensity changes of the
200 decodee rather than the weight distribution of the decoder. Moreover, the discordant performance between NDP
201 (nonmonotonic, **Fig 3A, C, and E**) and dot product (almost monotonic, **Fig 3B, D, and F**) suggests that previously
202 reported results ⁵ were primarily due to an increase in the magnitude of brain activity in specific regions, but in a
203 way that becomes less similar to the decoder. Yet, both NDP and dot product were robust to the removal of voxel
204 weights. These results again refute superiority of “optimal” decoders above that of a meta-contrast map decoder,
205 which was derived from GLM results. This reinforces the notion that location-only performs sufficiently, and that

206 useful information is provided only by the decodee activity within the locations where a decoder has non-zero
 207 weights.
 208
 209



210
 211 **Figure 2. Discrimination performance is similar for all three pain decoders and is a function of voxel locations, not weighted patterns.** (A)
 212 Meta-analysis of across-subject discrimination performance (AUC, chance = 0.5) for decoding pain from non-pain mental states for each of the
 213 three decoders, tested only for datasets independent of decoder derivation. On average, all decoders perform similarly. Square sizes indicate
 214 meta-analytic weight. (B-C) Across-subject decoding of pain from touch. (B) performance does not change when decoder pattern weights are
 215 distorted with increasing-size spatial smoothing. Sign = sign of voxel weights with 0 filtering, rendering decoder voxel values to only 0, -1, +1;
 216 filtering $\sigma = 0-20$ mm; ∞ = location-only. (C) Decoder performance depends only on a very small number of voxels, indicating information
 217 redundancy. The number of voxels constituting each decoder was systematically increased (from 10 voxels to the full decoder) and performance
 218 assessed for random samples of each size. 10% of each full decoder's voxel count (black ticks) discriminates pain from touch equivalently to the
 219 full decoders. Shades are standard deviations for spatial uncertainty, ignoring across-subject uncertainty. (D) Decoders were constructed using
 220 10% of the voxels in the full decoders, with voxels selected in order of their absolute magnitude (see Fig S7). The voxels with the highest
 221 absolute weights do not necessarily discriminate better than voxels with lower magnitudes, with the exception of pNsy in this dataset. (E)
 222 Selecting voxels based on their anatomical locations revealed that single regions (e.g., L. insula/parietal lobe) can discriminate similarly to the
 223 full decoders. Bars and shades are the 95% confidence intervals [CI] of means, except in C, where shades indicate standard deviations associated
 224 with permutation variability. In D and E, colored bars indicate the AUC of the full decoders.
 225



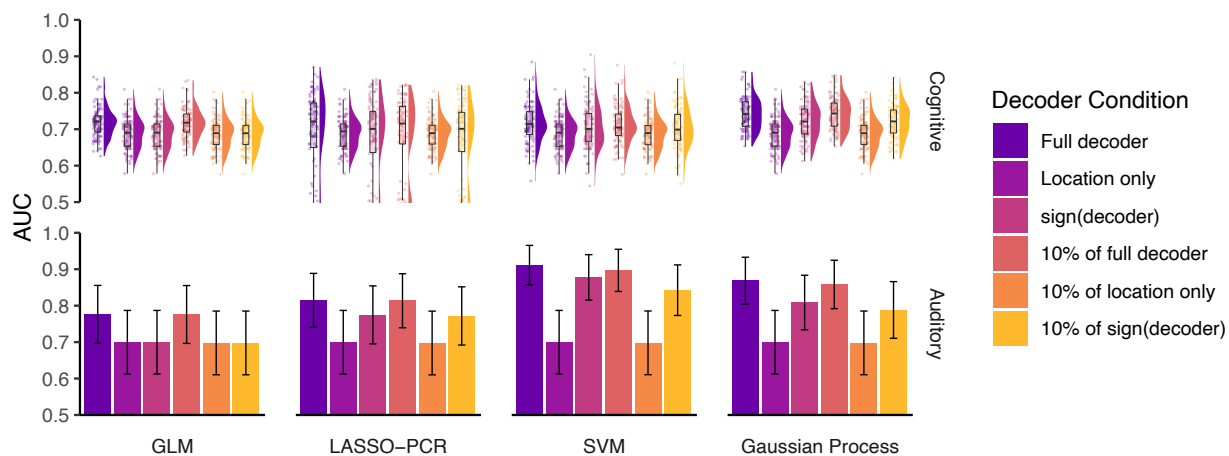
226
227
228
229
230
231

Figure 3. All decoders perform stimulus-perception mapping similarly, both with and without voxel weights. All three pain decoders perform equivalently, when location-only decoders are compared to the unfiltered-decoders, in mapping pain and heat perception ratings (A–B), mapping painful stimuli (C–D), and discriminating between pairs of painful stimuli (E–F). Nonmonotonic relationships indicate that the decoders cannot reliably predict subjective ratings or stimulus intensity. Vertical lines in A and B indicate the transition from heat (< 100) to pain (> 100). The dot products in B, D, and F were z-scored within each decoder for presentation purposes.

232 **Cognitive and Auditory Decoders Are Similarly Highly Redundant**

233 So far, we have shown that popular “optimized” pain decoders, as well as a meta-contrast map used as a decoder, are
 234 able to maintain their full performance after being perturbed and degraded, indicating that much of the information
 235 contained within them is superfluous. One worries that the findings may be specific to the modality studied, as pain
 236 and nociception are sensory systems for which no dedicated neocortical tissue has been uncovered in the cortex ²¹.
 237 As a result, there is long-standing debate as to specific or distributed encoding of pain perception (e.g., ²²; cf. ²³). To
 238 broaden our findings, we examined whether the uncovered principles apply to decoding for audition and language.
 239 Primary and secondary auditory cortex ^{24,25} are in close proximity to the somatosensory regions examined above for
 240 pain, while language representation with dedicated and functionally specific tissue is unique to humans ²⁶. We used
 241 data for language ¹⁹ and auditory ¹⁸ studies to construct decoders using task-specific contrast maps, SVM, LASSO-
 242 PCR, and Gaussian processes (our contrast maps closely resemble those reported in the original studies, **Fig S10–**
 243 **S11**; see Supplementary Methods). Our findings are entirely concordant with those for the pain decoders, in that all
 244 of the constructed decoders show similar performance, which was maintained after extreme perturbations (e.g., sign
 245 or location-only) (**Fig 4**), with only a few exceptions (see **Fig 4** comments). These findings generalize and provide
 246 compelling support for our main result: “optimized” decoders are highly redundant, and decoding primarily exploits
 247 information contained within voxel locations, independent of voxel weights. Moreover, task-specific GLM contrast
 248 maps are sufficient, implying that the meta-contrast maps are also not necessary.

249
250



251
252 **Figure 4. Different implementations of cognitive and auditory decoders perform similarly regarding discrimination**
 253 **performance and are robust to perturbations.**

254 We constructed decoders using general linear modeling (GLM), least absolute shrinkage and selection operating with principal
 255 components regression (LASSO-PCR), support vector machines (SVM), and Gaussian processes to decode **(top)** cognitive ¹⁹ and
 256 **(bottom)** auditory tasks ¹⁸. Much like the pain decoders, these decoders performed similarly and better than chance (chance = 0.5
 257 in both), and were relatively robust to perturbations. Just 10% of each decoder was enough to capture its full performance, and
 258 even extreme perturbations, such as 10% of the binary decoder or 10% of signed decoder, had little effect on performance. Error
 259 bars are the 95% confidence intervals of the AUCs. *Nota bene*, in the auditory task, discrimination performance is better with
 260 SVM and Gaussian Process than with GLM or LASSO-PCR. We suspect these differences are a consequence of specific
 261 instantiations of overfitting. We observed similar decoder-dependent performance variations for the pain decoders as well (see
 262 **Fig 2A**); yet, in further analyses none showed superiority over the others. In the auditory task, and for both SVM and Gaussian
 263 Process decoders, we also observed appreciable performance decrement for location-only and for 10% location-only decoders.
 264 This too was observed in the pain decoders. Like with the pain decoders, here, we also observed that sign-only decoders and 10%
 265 sign-only decoders performed similarly to the full decoders, again suggesting that negative weights at large scales can influence
 266 decoder performance.

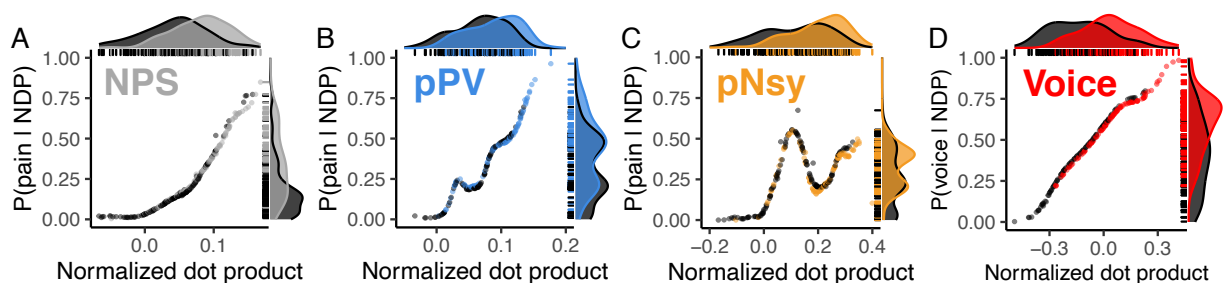
267

268 **Identification Remains A Challenge**

269 The ability of “optimized” decoders to *identify* mental states is repeatedly asserted in the literature^{3,5-11}, but
270 to our knowledge, remains untested. Specifically, the sensitivity and specificity reported in previous works are
271 estimated by changing decoder response thresholds for different pain stimuli⁵. If “optimized” decoders are used
272 with the objective of identification, then they should be able to pinpoint the specific mental state solely from the
273 similarity between the decoder and decodee, and, crucially, in the absence of a comparator and without such a
274 threshold tuning. In other words, identification should be based on a single observation and what we (or the decoder)
275 “know(s)” about the world. Therefore, instead of AUC, which implies a comparator, we tested identifiability by
276 calculating distributional overlap between the states of interest and no interest. Distributional overlap provides the
277 range of equal probability of belonging to the state of interest and state or states of no interest; here, equiprobability
278 implies unidentifiability. In addition, we were interested in assessing performance at the individual level. To do so,
279 we calculated the probability of a subject being in a specific mental state given that subject’s brain activity map. We
280 thus calculated distributional overlaps and state probabilities to assess the ability of decoders to identify mental
281 states.

282 Identification of pain states was similarly poor across the three pain decoders explored: overlaps between
283 states of interest and states of no interest were high ($\geq 68\%$) and the probabilities of being in pain (when actually in
284 pain) were low (median posterior probability ≤ 0.5) (**Fig 5a–c**). These results paint a markedly different picture than
285 the discrimination results, which simply show that NDPs *tend* to be greater when individuals are in pain; evidently,
286 adequate discrimination does not translate to identification. We built upon these findings by using the task-specific
287 contrast map decoder to decode audition of vocal versus non-vocal sounds¹⁸. While the performance of the voice
288 decoder was better than that of the pain decoders (overlap = 54%), it was still inadequate, as over half of the data
289 was unidentifiable (**Fig 5d**). The slight superiority of the voice decoder relative to the pain decoders may have
290 several explanations, including the homogeneity of the training and test sets used for the voice data, or simply that
291 some tasks are easier to identify than others. In any case, regardless of the mental state tested, identification
292 remained unreliable and thus is currently not feasible with fixed-weight decoders.

293



294 **Figure 5. Identification of mental states shows poor predictability.** Three pain decoders (NPS, pPV, and pNsy in A–C) and a voice decoder
295 (D) were used to test identification for mental states. x-axes are the normalized dot products between decoder and decodee, while y-axes are the
296 posterior probability of being in pain (A–C) or listening to voices (D). Distributions of normalized dot products and posterior probabilities
297 include both the decodee (light grey & colors) and comparator (dark grey) tasks. (A–C) Normalized dot products of the pain condition span the
298 entire distribution of comparator normalized dot products, and as a result, pain is not strongly isolated from the comparator conditions.
299 Quantitatively, this is evidenced by the strong decodee-comparator overlap for (A) NPS (overlap (95%CI) = 68% (59–82)), (B) pPV (79% (73–
300 90)), and (C) pNsy (73% (66–84)). This is reflected in the Bayesian model, which shows similar probabilities of being in pain for both pain and
301 pain-free conditions (each dot/line). To this end, all three decoders perform similarly, and cannot unequivocally *identify* pain, as indicated by their
302 low sensitivity/specificity (when specificity/sensitivity=0.95) of (NPS, A) 0.19/0.25, (pPV, B) 0.25/0.17, and (pNsy, C) 0.17/0.27. (D) In contrast
303 to pain, a contrast map decoder for identifying when a participant is listening to human voices separates more clearly the normalized dot products
304 of the decodee (red) from comparator (dark grey), but still performs poorly (overlap = 54% (46–66)). This separation is reflected in the Bayesian
305 model, which shows high probabilities when individuals are listening to human voices and lower probabilities when they are not. To identify the
306 mental state of listening to voices based on NDP with a specificity of 0.95, one would have a sensitivity of 0.19. Conversely, to identify the same
307 mental state with a sensitivity of 0.95, one would have a specificity of 0.48.

308

309

310

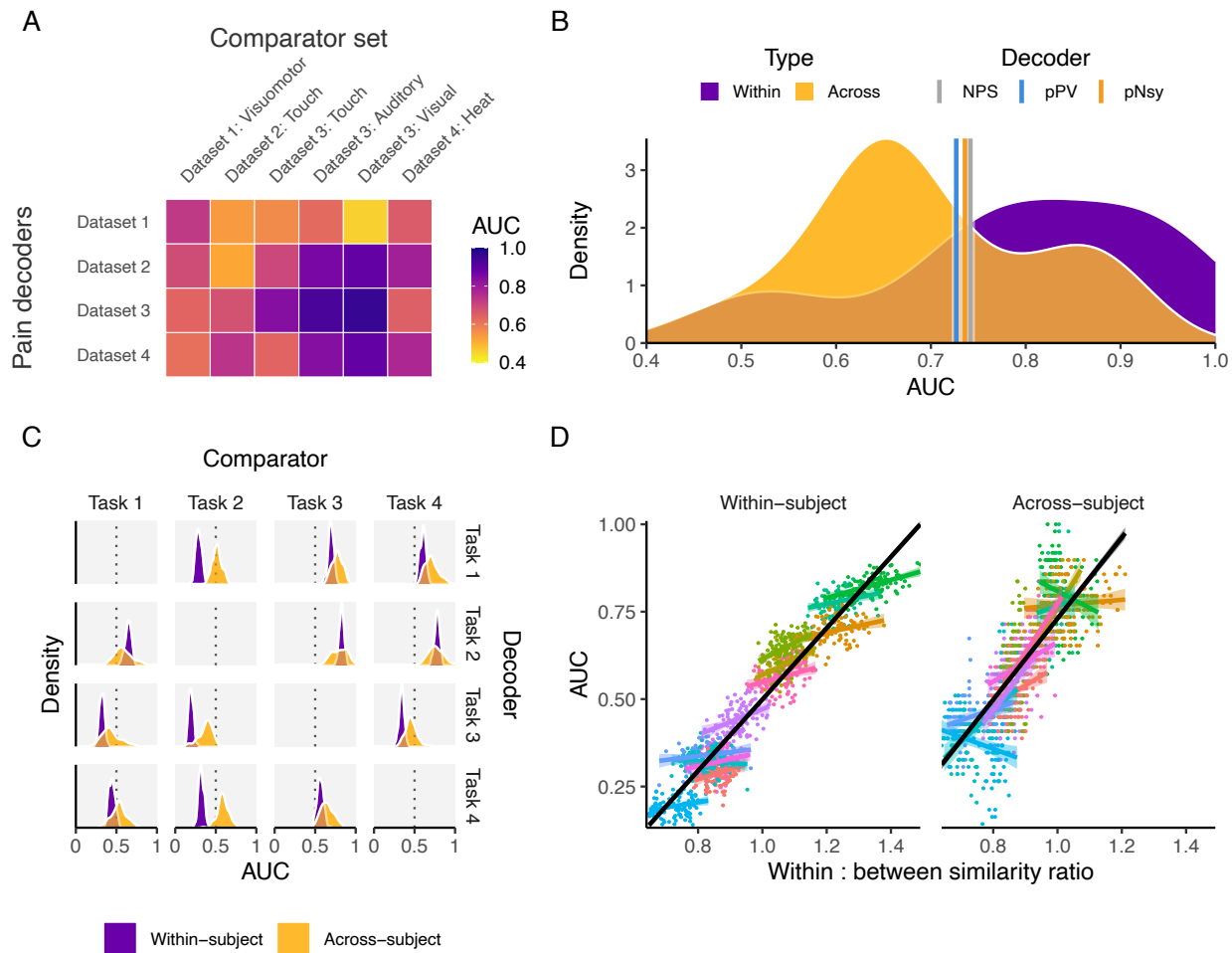
311 ***Brain Activity Maps Are Sufficient for Discrimination***

312 The similarity in performance achieved by meta-contrast maps or task-specific contrast maps and “optimized”
313 decoders prompted us to take another step back in the decoding derivation process. Given that pNsy is a composite
314 of contrasts from many studies (i.e., a meta-contrast pain decoder) and its decoding performance was similar to
315 “optimized” decoders (NPS and pPV), we assessed whether an *even* simpler construct—pain activity maps—was
316 sufficient to decode the state of being in pain. In other words, if no performance is lost by using contrast maps,
317 would task-derived GLM maps suffice as simpler but adequate decoders? Brain activity map decoders were created
318 using the average brain activity for each study’s pain task. Each activity map decoder was then used to discriminate
319 pain using the left-out brain activity maps of subjects both within and between studies (**Fig 6A**). Remarkably, these
320 decoders performed comparably to the ones presented hitherto (NPS, pPV, and pNsy), with an average within-study
321 AUC of 0.79 and between-study AUC of 0.69 (cf. ~0.73 for the fixed-weight decoders; **Fig 6B**). The lack of clear
322 superiority of “optimized” decoders relative to a meta-contrast map, and even simple activity map decoders, casts
323 serious doubt on the predictive and epistemological value of the more complex “optimized” decoders. These results
324 also raise the salient question: If decoding can be approached in so many different ways, what actually determines
325 decodability?

326 While decoding is difficult, decodability itself is likely predictable, yet to our knowledge remains
327 unexplored. To build upon our breed metaphor, some dogs exhibit features that largely overlap with other dogs, such
328 as the stature, color, and flat-faced features of pugs and French Bulldogs. Similarly, the mental state of “being in
329 pain” shares many features with other states; for example, unpleasantness, behavioral relevance, and saliency²⁷.
330 Therefore, the primary challenge of decoding is to tease apart these overlapping features. For this reason, it seems
331 logical that the similarity of activity maps within and between the decoder, decodee, and comparator would
332 determine decoding performance. If the decoder is built from activity maps that are dissimilar, the resulting average
333 map would have a low signal-to-noise ratio; if the decodees or comparators are dissimilar, then we can expect a
334 greater variance in NDPs; and if the decodees and comparators are similar to one another, then they will have a lot
335 of overlap and be difficult to tease apart. This logic implicates the neuroanatomical and physiological assumptions
336 previously mentioned, as heterogeneity across individuals should decrease similarity, making the NDPs more
337 variable and thus more difficult to discern. Using similarity metrics that reflect these relationships, we attempted to
338 explain decodability.

339 Until now, we have primarily focused on decoding across- rather than within-subjects. Intuitively, it is
340 apparent that, for many of the reasons elaborated above, decoding mental states should be more successful within-
341 subjects compared to across-subjects, as has been formulated by others^{28,29}. However, no systematic analysis of this
342 notion has been performed using fixed-weight decoders. Therefore, we investigated this question using data well-
343 suited for the question: fMRI data collected from 14 subjects who completed four cognitive tasks, each with 12
344 replicates¹⁹. These repetitions enabled the comparison of decoder performance within- and across-subjects. As
345 expected, decoding performance is more precise (smaller variance) within-subject (**Fig 6C**), but, interestingly, not
346 necessarily better (greater average AUC). We investigated whether the ratio of decodee to decodee-comparator
347 similarity (or within:between) can be a possible natural metric of why some decoders are more efficacious than
348 others. Higher performing decoders showed greater within:between ratios than lower performing decoders (**Fig 6D**).
349 Similarly, decoder similarity—the average NDP of all pairwise combinations of a decoder’s constituent activity
350 maps, a measure of reliability—could also explain much of the decoder performance, and in support of our previous
351 conclusions, this relationship is largely unaffected by thresholding and binarizing the decoder (**Fig S12**). Further
352 exploration showed that decodability, especially within-subject, is strongly predicated on these similarity metrics
353 (**Fig S13–S14; Table S1**). Decodee similarity, together with decodee-comparator similarity, is strongly predictive of
354 discriminability, accounting for up to 95% of the variance in AUCs. Our similarity metrics almost entirely explain
355 within-subject decodability, but only about 68% of AUC variance in across-subject decoding. This result may speak
356 to the assumptions violated by across-subject decoders, in that a similarity score across-subjects is less interpretable
357 than one calculated within a single subject since variance (e.g., brain anatomy) may be converted to bias (making all
358 brains fit the same template) during image preprocessing and registration.

359



360
361 **Figure 6. Decoders constructed from activity maps perform similarly to pattern-based decoders, and are dependent on both decodee and**
362 **comparator properties.** (A) Performance of four activity map decoders, based on the across-subject averaging for pain tasks, to differentiate
363 pain from six other mental states. (B) Among the activity map decoders, within study performance is slightly higher but extensively overlaps with
364 across study performance. Meta-analytic estimates of performance for NPS, pPV, and pNsy are within 0.4 standard deviation from the average
365 performance of both within and across study activity map decoders. (C-D) Properties of activity map decoders are examined within and across
366 subjects as a function of a cognitive task¹⁹. (C) Decoders (rows) are built from four cognitive tasks, tested on remaining three (columns), in a
367 within subject and across subject design. Within subject performance is always more consistent (i.e. it has smaller variance) but not necessarily
368 greater than across subject. For example, the within subject performance is always superior to across subject when using task 2 as the decoder.
369 The inverse is true when task 2 is the comparator, implying strong task dependence. (D) Decoder performance linearly scales with the ratio of
370 decodee similarity to decodee-comparator similarity (based on normalized dot product), for within- and across-subject comparisons. Because
371 discriminability depends on this ratio of similarities, they can be viewed as rules for decoding. Each color in (D) represents a decodee-comparator
372 pair of tasks 1-4 in (C); each point is a permuted sample; each colored line is the regression within a decodee-comparator pair; and the black line
373 is the regression across decodee-comparator pairs.

374
375 **Discussion**

376 In this study, we asked what the determinants, and limits, of decoding mental states are. We primarily emphasized
377 decoding pain, as this is the modality where the most emphatic claims have been made and where the “optimized”
378 decoders seem to have become accepted as enabling “mind-reading”³. For pain, audition, and language tasks, the
379 locations of a small subset of GLM-derived voxels were sufficient for achieving a discrimination of $AUC \approx 75\%$,
380 and a long list of machine learning tools could not consistently improve upon this performance. We also showed
381 that, in contrast to discriminating between states, identification of a given perceptual state is much harder. For the
382 first time, we advanced the concept of quantifying discriminability using a simple similarity metric, the NDP, with

383 which we provide models for within- and across-subject discrimination. The latter analyses indicated that
384 discriminability depends not only on the decoder, but also on similarity between the decodee and comparator.
385 Finally, we showed that, even in an example where within-subject discrimination was almost fully modeled with
386 similarity properties, there was a considerable decrease in the variance of across-subject discrimination that could be
387 explained. In doing so, we established limits of decodability based on the most popular linear models currently used
388 in fMRI literature.

389 Limitations of across-subject decoding and reverse inference have been acknowledged by others. For
390 example, the latest evidence shows that brain-behavioral phenotype associations seem to become reproducible only
391 with sample sizes of $N \gtrsim 2,000$ ³⁰. Yet, the extent of these limitations has not previously been quantified, nor has
392 decodability been modeled. Multiple approaches have been initiated to overcome these limitations. The simplest is
393 to constrain functional studies to within-subject investigations, thus bypassing the idiosyncrasies of anatomically
394 aligned, group-averaged results, but this approach also obviates across-subject decoding. The approach has been
395 used in various topics, including subject-specific localizers in vision³¹ and language studies³².

396 Perhaps the most widespread method is the multivoxel pattern analysis (MVPA). MVPA looks for
397 statistical evidence for *information* contained within groups of voxels (functional/anatomical regions of interest or
398 searchlights)³³. MVPA has been successful in decoding diverse brain states from fMRI activity patterns; for
399 example^{4,34-40}. MVPA typically uses subject-specific classifier models, and as a result, its accuracy drops when
400 predicting other subjects' responses^{28,29,38}. To extend MVPA results to across-subject applications, and to improve
401 on anatomical alignment, Haxby and colleagues^{33,41,42} developed an across-subject, high-dimensional, functional
402 alignment technique, named hyperalignment. It has previously been shown that hyperalignment, coupled with
403 MVPA, improves across-subject response classification to levels that are comparable to, or even better than, those
404 seen for within-subjects^{29,38}. To do so, hyperalignment exploits the temporal variability of stimulus-evoked brain
405 activity, yet it is designed to enable alignments based on diverse brain signals^{29,33,39}. Therefore, we explored
406 whether hyperaligning GLM brain activity maps would enhance across-subject decodability. In contrast to previous
407 work, our preliminary results did not show improvement between hyperaligned and anatomically aligned across-
408 subject decodability (data not shown). Still, it is possible that variants of hyperalignment may be useful in brain
409 activity-based across-subject decodability (e.g.,⁴³).

410 Our principal finding is consistent with the MVPA literature. A recent across-subject study used MVPA
411 (without hyperalignment) to uncover circuitry associated with pain relief commonly seen for eight different types of
412 analgesics⁴⁴. Their study identified brain locations involved in analgesia for multiple drugs, and thus it is consistent
413 with our main conclusion: location and not fixed patterns are sufficient decoders. In fact, in general, MVPA
414 identifies within- or across-subject brain locations, at macro- (voxel level) or micro-scale (sub-voxel level,⁴⁵), where
415 a certain discrimination or calculation is possible. The level of discrimination will ultimately be constrained by
416 differences in functional-anatomical coupling across individuals, in turn leading to distinct results within- and
417 across-subjects⁴⁶. Thus, specific underlying patterns may not be identified, which again is consistent and
418 complementary to our main conclusion.

419 Beyond promoting reverse inference, fixed-pattern decoders, also labeled as “signatures” or
420 “neuromarkers”^{2,3}, are purported to 1) unravel neural encoding of psychological constructs, 2) improve
421 decodability, 3) enable validation across studies and labs, and 4) provide falsifiable models. Since such decoders do
422 not outperform brain activity map-based decoders, we contend that the aforementioned assertions are untenable. It
423 follows that fixed-pattern decoders do not provide a defensible path forward for constructing a brain activity-derived
424 ontology of mental constructs⁴⁷.

425 Our demonstration that overlaying linear machine learning optimization algorithms did not improve on
426 linear contrast-derived decoders is not surprising. Indeed, similar conclusions have been reached in other domains,
427 such as medicine^{48,49}. Moreover, our findings support the idea that neuroimaging has not saturated the performance
428 of simple linear models⁵⁰. The reasons for this are manifold, and from a modeling viewpoint, it has been argued that
429 the added value of linear “machine learning” techniques is often small, exaggerated, and does not translate into
430 practical advantages⁵¹. Although unsurprising given the aforementioned work in this area, the apparent stark
431 discrepancy between our findings and those in the literature warrants explicit explanation.

432 How do we explain the discrepancy between our results and the literature, even when the same decoder is
433 used on the same data ⁵? We cannot escape the conclusion that “optimized” decoders are superfluous models.
434 Indeed, Wager and colleagues have also observed similar performance across several pain decoders, including NPS
435 and pNsy ^{5,52}. Moreover, the use of arbitrary performance metrics (here we base all analysis on NDP), lack of a
436 control (comparison to GLM modeling), and commonly mixing within- and across-subject performance metrics all
437 seem to mislead and propagate grandiose assertions ⁵³. In stark contrast, here we show that linear models—contrast
438 and activity maps—are capable of maximizing prediction, while being readily available and maintaining
439 interpretability. Yet, across-subject decodability remains complex; only brain location adds value, and depends on
440 within and between similarity of decoder, decodee, and comparator. These findings advance the general principles
441 of decoding mental states. Importantly, the limited and inadequate performance of fixed-weight across-subject
442 decoders, especially regarding identification, pose strict bounds on their utility in the domains of medical and legal
443 decision-making.

444

445 **Acknowledgments**

446 We would like to thank Dr. Thorsten Kahnt and Apkarian lab members for providing their thoughtful feedback.

447

448 **Authorship Contributions**

449 RJ, ADV, MNB, GDI, and AVA conceived the idea; RJ, ADV, JB, and LH performed the analyses; RJ, ADV, and
450 AVA drafted the manuscript; RJ, ADV, JB, GDI, and AVA edited the manuscript; and all authors approved the final
451 manuscript.

452

453 **Funding**

454 This work is funded by the National Institutes of Health (1P50DA044121-01A1). This material is based upon work
455 supported by the National Science Foundation Graduate Research Fellowship under Grant No. DGE-1324585.

456 **References**

- 457 1 Haynes, J. D. *et al.* Reading hidden intentions in the human brain. *Curr Biol* **17**, 323-328,
458 doi:10.1016/j.cub.2006.11.072 (2007).
- 459 2 Gabrieli, J. D., Ghosh, S. S. & Whitfield-Gabrieli, S. Prediction as a humanitarian and pragmatic contribution
460 from human cognitive neuroscience. *Neuron* **85**, 11-26, doi:10.1016/j.neuron.2014.10.047 (2015).
- 461 3 Kragel, P. A., Koban, L., Barrett, L. F. & Wager, T. D. Representation, Pattern Information, and Brain
462 Signatures: From Neurons to Neuroimaging. *Neuron* **99**, 257-273, doi:10.1016/j.neuron.2018.06.009 (2018).
- 463 4 Liang, M., Su, Q., Mouraux, A. & Iannetti, G. D. Spatial Patterns of Brain Activity Preferentially Reflecting
464 Transient Pain and Stimulus Intensity. *Cereb Cortex* **29**, 2211-2227, doi:10.1093/cercor/bhz026 (2019).
- 465 5 Wager, T. D. *et al.* An fMRI-based neurologic signature of physical pain. *N Engl J Med* **368**, 1388-1397,
466 doi:10.1056/NEJMoa1204471 (2013).
- 467 6 Wager, T. D. *et al.* A Bayesian model of category-specific emotional brain responses. *PLoS Comput Biol* **11**,
468 e1004066, doi:10.1371/journal.pcbi.1004066 (2015).
- 469 7 Poldrack, R. A., Halchenko, Y. O. & Hanson, S. J. Decoding the large-scale structure of brain function by
470 classifying mental States across individuals. *Psychol Sci* **20**, 1364-1372, doi:10.1111/j.1467-
471 9280.2009.02460.x (2009).
- 472 8 Woo, C. W., Roy, M., Buhle, J. T. & Wager, T. D. Distinct brain systems mediate the effects of nociceptive
473 input and self-regulation on pain. *PLoS Biol.* **13**, e1002036, doi:10.1371/journal.pbio.1002036 (2015).
- 474 9 Eisenbarth, H., Chang, L. J. & Wager, T. D. Multivariate Brain Prediction of Heart Rate and Skin
475 Conductance Responses to Social Threat. *J Neurosci* **36**, 11987-11998, doi:10.1523/JNEUROSCI.3672-
476 15.2016 (2016).
- 477 10 Marquand, A. *et al.* Quantitative prediction of subjective pain intensity from whole-brain fMRI data using
478 Gaussian processes. *Neuroimage* **49**, 2178-2189, doi:10.1016/j.neuroimage.2009.10.072 (2010).
- 479 11 Lindquist, M. A. *et al.* Group-regularized individual prediction: theory and application to pain. *Neuroimage*
480 **145**, 274-287, doi:10.1016/j.neuroimage.2015.10.074 (2017).
- 481 12 Kandel, E. R. *Principles of neural science*. 5th edn, (McGraw-Hill, 2013).
- 482 13 Gazzaniga, M. S. *The new cognitive neurosciences*. (MIT Press, 2000).
- 483 14 Feilong, M., Nastase, S. A., Guntupalli, J. S. & Haxby, J. V. Reliable individual differences in fine-grained
484 cortical functional architecture. *Neuroimage* **183**, 375-386, doi:10.1016/j.neuroimage.2018.08.029 (2018).
- 485 15 Mecacci, G. & Haselager, P. Identifying Criteria for the Evaluation of the Implications of Brain Reading for
486 Mental Privacy. *Sci Eng Ethics* **25**, 443-461, doi:10.1007/s11948-017-0003-3 (2019).
- 487 16 Chalmers, D. J. *The Conscious Mind: In Search of a Fundamental Theory*. (Oxford University Press, 1997).
- 488 17 Baliki, M. N., Geha, P. Y. & Apkarian, A. V. Parsing Pain Perception Between Nociceptive Representation
489 and Magnitude Estimation. *J Neurophysiol* **101**, 875-887, doi:10.1152/jn.91100.2008 (2009).
- 490 18 Pernet, C. R. *et al.* The human voice areas: Spatial organization and inter-individual variability in temporal
491 and extra-temporal cortices. *Neuroimage* **119**, 164-174, doi:10.1016/j.neuroimage.2015.06.050 (2015).
- 492 19 Jimura, K., Cazalis, F., Stover, E. R. & Poldrack, R. A. The neural basis of task switching changes with skill
493 acquisition. *Front Hum Neurosci* **8**, 339, doi:10.3389/fnhum.2014.00339 (2014).
- 494 20 Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C. & Wager, T. D. Large-scale automated
495 synthesis of human functional neuroimaging data. *Nat Methods* **8**, 665-670, doi:10.1038/nmeth.1635 (2011).
- 496 21 Chen, L. M. Cortical Representation of Pain and Touch: Evidence from Combined Functional Neuroimaging
497 and Electrophysiology in Non-human Primates. *Neurosci Bull* **34**, 165-177, doi:10.1007/s12264-017-0133-
498 2 (2018).
- 499 22 Segerdahl, A. R., Mezue, M., Okell, T. W., Farrar, J. T. & Tracey, I. The dorsal posterior insula subserves a
500 fundamental role in human pain. *Nat Neurosci* **18**, 499-500, doi:10.1038/nn.3969 (2015).
- 501 23 Iannetti, G. D. & Mouraux, A. From the neuromatrix to the pain matrix (and back). *Exp*
502 *Brain Res* **205**, 1-12, doi:10.1007/s00221-010-2340-1 (2010).
- 503 24 Brewer, A. A. & Barton, B. Maps of the Auditory Cortex. *Annu Rev Neurosci* **39**, 385-407,
504 doi:10.1146/annurev-neuro-070815-014045 (2016).
- 505 25 Fruhholz, S. & Grandjean, D. Multiple subregions in superior temporal cortex are differentially sensitive to
506 vocal expressions: a quantitative meta-analysis. *Neurosci Biobehav Rev* **37**, 24-35,
507 doi:10.1016/j.neubiorev.2012.11.002 (2013).
- 508 26 Broca, P. P. Perte de la Parole, Ramollissement Chronique et Destruction
509 Partielle du Lobe Antérieur Gauche du Cerveau. *Bulletin de la Société Anthropologique* **2**, 235 (1861).

- 510 27 Mouraux, A. & Iannetti, G. D. The search for pain biomarkers in the human brain. *Brain* **141**, 3290-3307,
511 doi:10.1093/brain/awy281 (2018).
- 512 28 Cox, D. D. & Savoy, R. L. Functional magnetic resonance imaging (fMRI) "brain reading": detecting and
513 classifying distributed patterns of fMRI activity in human visual cortex. *Neuroimage* **19**, 261-270,
514 doi:10.1016/s1053-8119(03)00049-1 (2003).
- 515 29 Haxby, J. V. *et al.* A common, high-dimensional model of the representational space in human ventral
516 temporal cortex. *Neuron* **72**, 404-416, doi:10.1016/j.neuron.2011.08.026 (2011).
- 517 30 Marek, S. *et al.* Towards Reproducible Brain-Wide Association Studies. *bioRxiv*, 2020.2008.2021.257758,
518 doi:10.1101/2020.08.21.257758 (2020).
- 519 31 Nasr, S., Polimeni, J. R. & Tootell, R. B. Interdigitated Color- and Disparity-Selective Columns within
520 Human Visual Cortical Areas V2 and V3. *J Neurosci* **36**, 1841-1857, doi:10.1523/JNEUROSCI.3518-
521 15.2016 (2016).
- 522 32 Fedorenko, E. & Blank, I. A. Broca's Area Is Not a Natural Kind. *Trends Cogn Sci* **24**, 270-284,
523 doi:10.1016/j.tics.2020.01.001 (2020).
- 524 33 Haxby, J. V., Connolly, A. C. & Guntupalli, J. S. Decoding neural representational spaces using multivariate
525 pattern analysis. *Annu Rev Neurosci* **37**, 435-456, doi:10.1146/annurev-neuro-062012-170325 (2014).
- 526 34 Kahnt, T. A decade of decoding reward-related fMRI signals and where we go from here. *Neuroimage* **180**,
527 324-333, doi:10.1016/j.neuroimage.2017.03.067 (2018).
- 528 35 Haynes, J. D. & Rees, G. Predicting the orientation of invisible stimuli from activity in human primary visual
529 cortex. *Nat Neurosci* **8**, 686-691, doi:10.1038/nn1445 (2005).
- 530 36 Pilgramm, S. *et al.* Motor imagery of hand actions: Decoding the content of motor imagery from brain activity
531 in frontal and parietal motor areas. *Hum Brain Mapp* **37**, 81-93, doi:10.1002/hbm.23015 (2016).
- 532 37 Oosterhof, N. N., Tipper, S. P. & Downing, P. E. Viewpoint (in)dependence of action representations: an
533 MVPA study. *J Cogn Neurosci* **24**, 975-989, doi:10.1162/jocn_a_00195 (2012).
- 534 38 Al-Wasity, S., Vogt, S., Vuckovic, A. & Pollick, F. E. Hyperalignment of motor cortical areas based on motor
535 imagery during action observation. *Sci Rep* **10**, 5362, doi:10.1038/s41598-020-62071-2 (2020).
- 536 39 Haxby, J. V., Gobbini, M. I. & Nastase, S. A. Naturalistic stimuli reveal a dominant role for agentic action
537 in visual representation. *Neuroimage* **216**, 116561, doi:10.1016/j.neuroimage.2020.116561 (2020).
- 538 40 Kaplan, J. T. & Meyer, K. Multivariate pattern analysis reveals common neural patterns across individuals
539 during touch observation. *Neuroimage* **60**, 204-212, doi:10.1016/j.neuroimage.2011.12.059 (2012).
- 540 41 Haxby, J. V., Guntupalli, J. S., Nastase, S. A. & Feilong, M. Hyperalignment: Modeling shared information
541 encoded in idiosyncratic cortical topographies. *Elife* **9**, doi:10.7554/eLife.56601 (2020).
- 542 42 Guntupalli, J. S. *et al.* A Model of Representational Spaces in Human Cortex. *Cereb Cortex* **26**, 2919-2934,
543 doi:10.1093/cercor/bhw068 (2016).
- 544 43 McNorgan, C., Smith, G. J. & Edwards, E. S. Integrating functional connectivity and MVPA through a
545 multiple constraint network analysis. *Neuroimage* **208**, 116412, doi:10.1016/j.neuroimage.2019.116412
546 (2020).
- 547 44 Duff, E. P. *et al.* Learning to identify CNS drug action and efficacy using multistudy fMRI data. *Sci Transl*
548 *Med* **7**, 274ra216, doi:10.1126/scitranslmed.3008438 (2015).
- 549 45 Op de Beeck, H. P. Probing the mysterious underpinnings of multi-voxel fMRI analyses. *Neuroimage* **50**,
550 567-571, doi:10.1016/j.neuroimage.2009.12.072 (2010).
- 551 46 Clithero, J. A., Smith, D. V., Carter, R. M. & Huettel, S. A. Within- and cross-participant classifiers reveal
552 different neural coding of information. *Neuroimage* **56**, 699-708, doi:10.1016/j.neuroimage.2010.03.057
553 (2011).
- 554 47 Lenartowicz, A., Kalar, D. J., Congdon, E. & Poldrack, R. A. Towards an ontology of cognitive control. *Top*
555 *Cogn Sci* **2**, 678-692, doi:10.1111/j.1756-8765.2010.01100.x (2010).
- 556 48 Desai, R. J., Wang, S. V., Vaduganathan, M., Evers, T. & Schneeweiss, S. Comparison of Machine Learning
557 Methods With Traditional Models for Use of Administrative Claims With Electronic Medical Records to
558 Predict Heart Failure Outcomes. *JAMA Netw Open* **3**, e1918962, doi:10.1001/jamanetworkopen.2019.18962
559 (2020).
- 560 49 Christodoulou, E. *et al.* A systematic review shows no performance benefit of machine learning over logistic
561 regression for clinical prediction models. *J. Clin. Epidemiol.* **110**, 12-22, doi:10.1016/j.jclinepi.2019.02.004
562 (2019).
- 563 50 Schulz, M. A. *et al.* Different scaling of linear models and deep learning in UKBiobank brain images versus
564 machine-learning datasets. *Nat Commun* **11**, 4238, doi:10.1038/s41467-020-18037-z (2020).

- 565 51 Hand, D. J. Classifier Technology and the Illusion of Progress. *Statist. Sci.* **21**, 1-14,
566 doi:10.1214/088342306000000060 (2006).
- 567 52 Geuter, S. *et al.* Multiple Brain Networks Mediating Stimulus-Pain Relationships in Humans. *Cereb Cortex*
568 **30**, 4204-4219, doi:10.1093/cercor/bhaa048 (2020).
- 569 53 Hu, L. & Iannetti, G. D. Painful Issues in Pain Prediction. *Trends Neurosci* **39**, 212-220,
570 doi:10.1016/j.tins.2016.01.004 (2016).
571