1

Title**: Activity-dependent organization of prefrontal hub-networks for associative learning and signal transformation**

**Authors**

Masakazu Agetsuma[1, 2, 3, 4] *, Issei Sato[5], Yasuhiro R Tanaka[6], Luis Carrillo-Reid[7], Atsushi Kasai[8], Yoshiyuki Arai[3], Miki Yoshitomo[1], Takashi Inagaki[1], Hitoshi Hashimoto[8, 9, 10, 11, 12], Junichi Nabekura[1], and Takeharu Nagai[3]

**Affiliations**

1, Division of Homeostatic Development, National Institute for Physiological Sciences, 38 Nishigohnaka Myodaiji-cho, Okazaki, Aichi, 444-8585, Japan

2, Japan Science and Technology Agency, PRESTO, 4-1-8 Honcho, Kawaguchi, Saitama, 332-0012, Japan

3, SANKEN (The Institute of Scientific and Industrial Research), Osaka University, Mihogaoka 8-1, Ibaraki, Osaka 567-0047, Japan

4, Division of Molecular Design, Research Center for Systems Immunology, Medical Institute of Bioregulation, Kyushu University, 3-1-1 Maidashi, Higashi-ku, Fukuoka 812-8582, Japan

5, Department of Computer Science, Graduate School of Information Science and Technology, The University of Tokyo. 7-3-1 Hongo, Bunkyo-ku, Tokyo 113-0033, Japan.

6, Brain Science Institute, Tamagawa University, 6-1-1 Tamagawagakuen, Machida, Tokyo, 194-8610, Japan

7, Instituto de Neurobiologia, National Autonomous University of Mexico, Boulevard Juriquilla 3001, Juriquilla, Queretaro, CP 76230, Mexico

8, Graduate School of Pharmaceutical Sciences, Osaka University, Yamadaoka 1-6, Suita, Osaka 565-0871, Japan

9, United Graduate School of Child Development, Osaka University, Kanazawa University, Hamamatsu University School of Medicine, Chiba University, and University of Fukui, 2-2 Yamadaoka, Suita, Osaka 565-0871, Japan

10, Division of Bioscience, Institute for Datability Science, Osaka University, 1-8 Yamadaoka, Suita, Osaka 565-0871, Japan

11, Open and Transdisciplinary Research Initiatives, Osaka University, 2-1 Yamadaoka, Suita, Osaka 565-0871, Japan

12, Graduate School of Medicine, Osaka University, 2-2 Yamadaoka, Suita, Osaka 565-0871, JapanUse superscript numbers (1, 2, 3) to designate author affiliations.

*Corresponding author. Email: age@nips.ac.jp

**Abstract**

Associative learning is crucial for adapting to environmental changes. The encoding of associative learning involves the dorso-medial prefrontal cortex (dmPFC), and is underpinned by interactions within the resident neuronal population. However, the nature of this population coding is poorly understood. Here we developed a pipeline for computational dissection and longitudinal two-photon imaging of neural population activities in the mouse dmPFC during fear-conditioning procedures, enabling us to detect learning-dependent changes in the dmPFC topology. Through regularized regression methods and graphical modeling, we found fear conditioning organized neuronal ensembles encoding conditioned responses (CR), with enhancing their coactivity, functional connectivity, and association with conditioned stimuli (CS). This suggests that fear conditioning drives dmPFC reorganization to generate novel associative circuits for CS-to-CR transformation. Importantly, neurons strongly responding to unconditioned stimuli (US) during conditioning anterogradely became a hub of the CR ensemble. Altogether, we demonstrate learning-dependent dynamic modulation of population coding structured on an activity-dependent hub-network formation within the dmPFC.

**Teaser**

Optical and computational dissection uncovered how prefrontal cortical networks are rewired to encode new associative memory

**Significance statement**

Animals learn to adapt to changing environments. Associative learning is one of the simplest types of learning that has been intensively studied over the past century. Recent development in molecular, genetic, and optogenetic methods has enabled the identification of a neural population encoding the associative memory in the brain. However, it remains unclear how information is stored and processed by the neural population to encode and retrieve the associative memory. To investigate the nature of this population coding, we developed an optical and computational dissection method, demonstrating how associative learning drives reorganization of the neural network in the dorso-medial prefrontal cortex and generates novel circuits for associative memory and signal transformation.

**MAIN TEXT**

**Introduction**

Animals learn to adapt to changing environments for survival. Associative learning, such as classical conditioning, is one of the simplest types of learning that has been intensively studied over the past century (*1, 2*). It is based on repeated pairings of a neutral conditioned stimulus (CS) such as a tone, and an unconditioned stimulus (US) such as foot shock, that eventually makes the subjects respond to the CS by itself and elicit a conditioned response (CR), e.g., freezing behavior in the associative fear learning paradigm. During the last two decades, technical development in molecular, genetic, and optogenetic methods has enabled the identification of a population of neurons in the brain, named the memory engram, which encodes and regulates associative memory (*3*). How information is stored and processed by the neural population to encode and retrieve the associative memory, however, remains unclear (*3*). In addition, although previous studies proposed the possibility that the formation of associative memory may involve novel associative connections between the originally distinct CS and US networks to enable the CS-to-CR transformation, direct evidence is quite limited.

The prefrontal cortex (PFC) is a brain region that regulates associative fear memory, which is evolutionarily conserved in mammals, from humans to primates to rodents (*4-9*). Dysfunction of the PFC may lead to various psychiatric diseases including the post-traumatic stress disorder (*10*), and the associative fear memory paradigm has been used as a research model to investigate the underlying mechanisms of this disorder. The dorsal part of the medial prefrontal cortex (dmPFC) of rodents is a brain region demonstrated to be important for the retrieval of associative fear memory (*11-16*). During fear memory retrieval and evoked freezing responses (i.e., CR), activated individual neurons (*17*) or an enhanced synchrony of neural populations (*14*) in the dmPFC are observed, while pharmacological or optogenetic silencing of the dmPFC and its projections to specific downstream targets suppresses fear memory retrieval (*11, 12*), revealing that associative fear memory is normally stored in the dmPFC. Therefore, the dmPFC can serve as an interesting target to address the fundamental question of what structural and computational alterations in the prefrontal networks are required to organize novel associative memories. Also, the study may contribute to our understanding of how novel associative memory is stored in the dmPFC together with pre-existing networks such as those regulating sensory and motor information.

To address these points, here we developed a pipeline for computational dissection and longitudinal imaging of neural population activities in the dmPFC during fear conditioning procedures in mice, which enabled us to uncover learning-dependent changes in the internal topology, functional connectivity and computational architecture of the dmPFC upon memory acquisition.

109

110

**Results**

112

**Longitudinal imaging of neural population activities in mouse dmPFC during fear-conditioning procedures**

To perform longitudinal imaging of neuronal population activities in the dmPFC during fear-conditioning procedures in mice, we first developed a system to perform cued-fear conditioning and memory retrieval while imaging awake and behaving mice with a two-photon microscope (Fig. 1), which enabled us to record the neural activities of hundreds of neurons with single-cell resolution and to further elucidate experience-dependent changes in functional connectivity in the dmPFC internal network, as shown later. The mice were head-fixed under the microscope objective and placed on a running disk, which was used to record the mouse locomotion status (e.g., locomoting, stationary, or expressing a freezing response) (Fig. 1A). Tones and foot shocks were delivered as the CS and US, respectively. Two different tones were used; one was associated with the US (CS+) and the other was not (CS−), as described in previous studies (*14, 18*). On day 3 (D3), the mice underwent a habituation session, in which they received 4 presentations of the CS− and CS+ alternately. The habituation session was immediately followed by discriminative fear conditioning session on the same day, in which only the CS+ was paired with the US (Fig. 1B). The US was delivered during the last 1-s of each 30-s CS+ trial. The CS− and CS+ trials were performed alternately (inter-trial intervals, 50–150 s). The next day (D4), the conditioned mice underwent a retrieval and extinction session, in which they received 4 presentations of the CS− and 12 presentations of the CS+ (4 presentations of the CS− and CS+ trials alternately, followed by 8 CS+ trials; Fig. 1B). Behavioral analyses revealed that the mice learned to exhibit freezing-like behavior, i.e., decrease their locomotion as a conditioned response (CR), specifically during the CS+, only after the fear conditioning (Fig. 1B, C). On the other hand, as reported previously (*14, 18*), the CR observed during the early phase on D4 was extinguished after repeated exposure to the CS+ only (Figs. 1D, E). Overall, these behavioral data established that our behavioral system and the fear conditioning protocol were useful for observing a change in the neural representation during associative fear learning.

Next, to monitor the neural activities in the dmPFC by two-photon microscopy, we implanted a 2-mm microprism along the rostral midline of the brain to optically access the dmPFC region. Although the size of the prism was larger than that of prisms used in previous work (*19*), there was sufficient space and no callosal fibers between the hemispheres around the dmPFC area,

143 especially at the rostral region, enabling the smooth insertion of the prism without cutting prefrontal

144 or callosal neural fibers (Fig. 2A). Using a genetically encoded $Ca^{2+}$ indicator, GCaMP6f, expressed

145 by an adeno-associated virus (AAV), the activities from a wide region of the prefrontal area were

146 chronically visualized (Fig. 2B, C and Movie S1). To specifically record the activity of the

147 excitatory neurons (20) and separate them from inhibitory neurons that may have a distinct function

148 in the dmPFC (13), the GCaMP6f was expressed under regulation of the CaMKII promoter (21,

149 22). The CS and US presentation did not disturb image acquisition (Movie S2). We focused on

150 analyzing activities in the dmPFC area (see the Materials and Methods for details). In most of the

151 data analyses, the neural representation during the first 3 trials of the fear conditioning on D3 (D3-

152 early [D3E]) were compared with those during the first 3 trials on D4 (D4-early [D4E]) to assess

153 the changes occurring after the fear conditioning and memory consolidation. The data obtained

154 during the last 3 trials on D3 (D3-late [D3L]) were used to assess the late conditioning phase, and

155 the data obtained during the last 3 trials on D4 (D4-late [D4L]) were used to assess the extinction

156 phase.

157 Prior to investigating population coding in the dmPFC, we first summarized the single-

158 neuron responses to the CS+ and CS− before and after acquisition of the fear memory (Figs. 2D-F

159 and S1). We found that approximately 60% of neurons exhibited a significant change in neural

160 activity during the CS+ and/or CS−, and approximately 20% of neurons showed significant

161 responses to both the CS+ and CS−. The distributions of these types of neurons were consistent

162 throughout the learning process (Figs. 2F and S1). This type of "mixed selectivity" (responsive to

163 variable task-relevant aspects) has been reported in the primate PFC (23) as well as in the mouse

164 caudal mPFC during a decision-making task (24). The potential advantage of the mixed selectivity

165 was proposed to enhance the number of tasks that each neural circuit, with a limited number of

166 neurons, can handle, through high-dimensional neural representations implemented by a population

167 of neurons (23, 25). This encouraged us to further analyze the population coding for fear memory.

168

169 **Newly emerged and unique neuronal ensembles in dmPFC encoding conditioned responses**

170 Our goal in this study was to dissect the computational architecture composed by a neural

171 population in the dmPFC enabling the distinctive acquisition of a novel associative memory. For

172 this purpose, we first extracted a group of neurons encoding the conditioned response (named CR

173 ensemble) (Fig. 2G, H), and compared it with the neurons encoding regular locomotion (RL) to

174 evaluate the uniqueness of the CR ensemble (Figs. 3 and S2). As methods to analyze neural

175 architecture embedded in the neural population activities, previous studies utilized unsupervised

176 algorithms such as Principal Component Analysis or Non-Negative Matrix Factorization (14, 26,

177   *27*). These algorithms first seek and distinguish embedded structures in neural data without

178   considering the behavioral labels (e.g., freezing responses), which are further used to test which

179   extracted structure is most likely to correlate with and explain respective behaviors, e.g., behavior

180   A and B. By such methods, the identified ensembles corresponding to behavior A and the identified

181   ensemble corresponding to behavior B may become dissimilar from each other as a result of

182   methodological bias, irrespective of the actual neural architectures. In the present study, instead of

183   these methods, we introduced a supervised and model-based decoding algorithm, named elastic net

184   (*28*) (Figs. 2G-H). The elastic net is a regularization and variable selection algorithm based on the

185   regression model (Fig. 2G; see the Materials and Methods for details) (*28*). This method enabled

186   us to independently extract the CR ensemble and RL ensemble from the same mice (Figs. S2), and

187   is thus helpful for further comparing the ensembles and systematically verifying whether neurons

188   in the CR ensembles were unique or mostly overlapped with the RL ensembles.

189          We extracted CR ensembles using neural activity data obtained during the CS+ presentation

190   of D4E (retrieval session), and evaluated the fitting and decoding performance of the obtained

191   model (Fig. 2H) after optimizing the hyper parameter "alpha" for the elastic net as explained below

192   (see the Materials and Methods for details). Compared with the conventional regularization

193   algorithm Lasso, the advantage of the elastic net is that this procedure enables us to optimize the

194   size of the selected population, especially when an analyzed neural network includes strongly

195   correlated neural pairs, which is likely the case for our data considering the results shown below

196   and previous electrophysiological observations (*14*). To carefully verify the overlap between the

197   CR ensemble and RL ensemble, it is important to avoid missing neurons encoding the respective

198   information. For this purpose, we evaluated remaining information encoded by neurons excluded

199   from the CR ensemble at each alpha, and defined optimal alpha as the one that minimize such

200   remaining information (Figs. S2C, D and S3; see also Materials and Methods). A wide range of the

201   alpha values for the CR ensemble of each individual mouse was tested (Fig. S3). This systematic

202   optimization procedure revealed a general trend that a larger alpha tended to select a smaller number

203   of CR ensemble neurons (Fig. S3B, top), and though the decoding performance of the smaller

204   number of selected CR ensembles was very high, equivalent to that of the others (Fig. S3B, middle),

205   the removal of such a smaller portion from the whole set of neurons was not sufficient to

206   substantially diminish the information encoded by the remaining neurons for some mice (Figs. S3B,

207   top and bottom, and S3D-F). This suggests that the CR was redundantly encoded in the dmPFC.

208   On the other hand, alpha values didn't clearly affect the discrimination of the RL ensemble (Fig.

209   S4). After determining the optimal alphas for individual circuits, we confirmed a substantial

210   reduction of the decodability by removing all the selected CR ensemble neurons (Figs. S2D, and

211  S3) and verified that a sufficiently large portion of the dmPFC neurons encoding the CR were

212  selected as the CR ensemble neurons.

213       We eventually confirmed that the CR ensemble obtained by the optimal alpha was highly

214  predictive for the CR during the retrieval session (see the example shown in Fig. 2H; mean ± SE of

215  the prediction accuracy, $0.9450 \pm 0.0265$, N=7 mice; individual data are shown later in Fig. 3G).

216  As for the spatial distribution, the identified CR ensemble neurons were spatially intermingled over

217  the field of view, as shown in Figs. 2H and 3B.

218       Following the optimization of the hyper parameter alpha, we evaluated the specificity and

219  uniqueness of the extracted CR ensemble. We confirmed that most of the neurons involved in the

220  CR ensemble were unique and did not overlap with the RL ensemble (Figs. 3A-C). We then

221  conceived the hypothesis that the unique CR ensemble might dominantly and exclusively explain

222  the behaviors of the mice during CS +-evoked memory retrieval as an encoder of the acquired

223  associative memory. If this is true, RL ensembles, distinct from CR ensembles (Fig. 3C), should

224  have diminished decodability for the behavior during CS+ during memory retrieval. To test this

225  possibility, we checked the decoding performance of the RL ensembles for the behaviors observed

226  during the CS+ at each of the learning steps (Figs. 3D-G and S5).  The decoding performance by

227  the RL ensemble to the RL was similar between pre- and post-memory consolidation (Fig. 3E). The

228  decoding performance of the RL ensemble to the behaviors during CS+ presentation at D3E (before

229  fear memory consolidation) was similar to that for the RL (Figs. 3D, F). In contrast, the decoding

230  performance of the RL ensemble to the behaviors during the CS+ on D4E (during fear memory

231  retrieval) was significantly reduced compared with that of D3E (Figs. 3D, F). There was a small,

232  but not significant, change during the fear conditioning (D3E vs D3L; Fig. S5), and importantly,

233  the reduced decodability of the behavior during CS+ at D4E (memory retrieval) was substantially

234  recovered after the extinction training (no significant difference between D3E and D4L, and a

235  significant difference between D4E and D4L; Figs. 3F and S5). On the other hand, the decodability

236  of CR ensembles was specific to the CR and not applicable to the RL on D4 (Fig. 3G). These results

237  established that the CR, or the behavior during the memory retrieval, was dominantly explained by

238  the CR ensembles, supporting the idea that the CR ensemble systematically extracted was a

239  dominant and specific group of neurons encoding the CR during memory retrieval, emerged after

240  consolidation of the fear memory and was suppressed by extinction.

241

242  **Coactivity within the CR ensemble was specifically enhanced after fear conditioning**

243       In these CR ensembles, we observed a slight but significant increase in CS+ activatable

244  neurons, but no change in CS+ inactivated neurons after fear conditioning (Fig. S6). In contrast,

245 other cells (neurons that were not included in the CR ensembles: Non-CR ensemble [Non-CRE]

246 neurons) exhibited no significant changes in the CS+ activatable neurons, with a significant increase

247 in CS+ inactivated neurons. Neurons in the RL ensembles did not exhibit any significant change in

248 CS+ responsiveness. We detected no significant change in CS– responsiveness in any of the

249 categories. Because these CR ensembles were discriminated by the data and behavioral labels

250 during the CS+, not by comparisons between those during the CS+ and those during the presentation

251 of other stimuli, our method produced no bias toward the CS+ during selection of the CR ensemble

252 neurons. These results indicated that there might be some mechanism that makes neurons involved

253 in the CR ensembles dominantly activated by the CS+ after memory consolidation.

254      To further analyze and characterize the identified CR ensemble toward elucidating the

255 mechanism underlying associative learning, we measured the change in the coactivity of the neural

256 network during the CS+ presentation by comparing the pairwise correlation coefficients (R) (*29*)

257 between pre- and post-memory consolidation. We found that, after the fear conditioning, only the

258 positively correlated fraction was enhanced specifically within the CR ensemble, and not in the

259 outside network (Non-CRE) (Fig. S7A). Statistical analyses demonstrated that this enhancement in

260 positive correlation after the fear conditioning, as well as the enhanced ratio of significantly and

261 positively correlated pairs, specifically occurred in the CR ensemble (Figs. S7A-C). Analyses based

262 on the shuffled data, where the activity of each neuron was preserved but the temporal order was

263 randomly shuffled neuron by neuron, revealed no significant difference between the CR ensemble

264 and Non-CRE (Figs. S7A, C), suggesting that the specific enhancement of the coactivity of the CR

265 ensemble in the real data did not derive from the enhanced neural activation. Similar enhancement

266 of the coactivity was observed in the CR ensemble excluding the RL-ensemble overlapped neurons

267 (Figs. S7A-C). In addition, changes in the coactivity across the categories (coactivity between CR

268 ensembles and Non-CRE) were significantly smaller than those within the CR ensembles (Fig. S7C).

269 These results led us to hypothesize that the functional connectivity within the CR ensemble was

270 specifically enhanced as a result of the fear conditioning, contributing to enhance the coactivity.

271

272 **Enhanced internal connectivity and association with conditioned stimuli (CS) in the CR**

273 **ensemble after fear conditioning**

274      To test the hypothesis above, we introduced a probabilistic graphical model method, the

275 conditional random field (CRF) model (*30, 31*). This method evaluates the conditional probability

276 that a group of neurons fire together given that one neuron is active (Fig. 4A). Among the various

277 mathematical algorithms used to evaluate possible functional connectivity of neural networks and

278 ensembles, the CRF model is one of the most reliable methods because the results of the calculation

279 (functional connectivity) have already been carefully evaluated by two-photon holographic
280 optogenetics and consequential behavioral modulation (*30, 31*).

281       Using this method, we found that, after the fear conditioning (D4E), the functional
282 connectivity was significantly higher in the CR ensemble (Fig. 4B). This method also allowed us
283 to evaluate the information coding of any arbitrary label, e.g., CS+ or CS−, and we found that the
284 CS+ information encoded by the CR ensemble was significantly higher than that of Non-CRE (Fig.
285 4C). Importantly, our method did not produce any bias to the CS+ in selecting CR ensemble neurons,
286 as explained above. Therefore, this result indicates that the neural population encoding the CR was
287 dominantly associated with the CS+ information. In addition, we found that the enhancement in
288 both the functional connectivity and CS+ predictability was experience-dependent and derived after
289 the fear conditioning, dominantly in the CR ensemble neurons (Figs. 4D, E). In contrast, the
290 changes in information coding for the CS− were not significantly different between the CR
291 ensemble and the Non-CRE (Fig. 4E). Therefore, the emergence of the CR ensemble after fear
292 conditioning was accompanied by the enhancement of the internal coactivity, functional
293 connectivity and association with the CS+ selectively within the CR ensemble neurons, indicating
294 that fear conditioning drives dmPFC reorganization to generate novel associative circuits for CS-
295 to-CR transformation.

296

**297 Neurons responding to US during fear conditioning anterogradely became a hub of the CR**
**298 ensemble**

299       Finally, we hypothesized that the dmPFC reorganization that we observed after fear
300 conditioning might occur via activity-dependent modulation during the repeated CS+-US pairing.
301 This led us to search for the signature of this plasticity.

302       During the fear conditioning, we observed that some of the dmPFC neurons strongly
303 responded to the US (Fig. 4F). Interestingly, statistical analyses demonstrated that neurons
304 responsive to the US during fear conditioning were predominantly and significantly more involved
305 in the CR ensemble after the fear conditioning (Figs. 4G, H). This suggests that these US-responsive
306 neurons (USR) were preferably integrated into the CR ensemble, in which functional connectivity
307 might also be modulated and strengthened by US-evoked activity, perhaps together with the paired
308 CS+ signal.

309       To test this possibility, we performed further analyses based on the CRF modeling. We found
310 that the USR became functionally more connected within the CR ensemble than non-US responsive
311 neurons, while these differences were not observed in Non-CRE (Fig. 4I). This higher connectivity
312 was a result of the fear conditioning (Fig. 4J). The information coding for the CS+ was also

313 significantly higher in the USR, specifically in the CR ensemble (Fig. 4K), suggesting that the US-
314 responsive network was dominantly associated with the CS+ network when it became integrated
315 into the newly emerged CR ensemble. According to a previous study, higher functional connectivity
316 and higher decoding performance of sensory stimuli are typical features of pattern completion cells
317 whose activation could efficiently enhance the entire ensemble activity for a specific sensory
318 stimulus and promote the stimulus-associated behaviors of mice (*30*). These results collectively
319 suggest that the USR in the dmPFC became a hub of the novel neural ensemble linking the CS+ to
320 the CR, a memory-evoked behavior, after the repeated CS+ and US parings.

321
322 **Discussion**
323 Altogether, our results based on the combination of methods for computational dissection and
324 longitudinal recording in the dmPFC demonstrate learning-dependent dynamic modulation of
325 population coding for associative fear learning, structured on an activity-dependent hub-network
326 formation within the dmPFC. Through regularized regression methods and graphical modeling, we
327 found that the repeated CS+-US pairing for the associative learning drives the dmPFC
328 reorganization to generate novel and unique neural circuits for CS-to-CR transformation, with
329 enhanced internal coactivity, connectivity, and association with the CS+. Upon this prefrontal
330 reorganization, neurons activated by the US during fear conditioning were anterogradely and
331 predominantly integrated into the CR ensemble. The eventual network stemming from these USR
332 gained typical features of pattern completion cells of the CR ensemble, which are supposed to work
333 as a hub in the prefrontal networks to predominantly relay the CS+ information and promote the
334 CR (Fig. S8).

335 To our knowledge, this is the first in vivo evidence directly demonstrating that the prefrontal
336 neural circuit for the associative memory was actually built based on an enhanced association
337 between the US network and the CS+ network as a result of CS+-US pairing and triggered network
338 reorganization. More than 60 years ago, Hebb proposed that repeated coactivation of a group of
339 neurons might create a memory trace through the enhancement of connections (*32*). Our results
340 suggest that Hebbian plasticity (i.e., fire together, wire together) might underlie the reorganization
341 of the prefrontal network structure during associative learning, enabling the emergence of a strong
342 link between the US signaling pathway and the CS+ signaling pathway to form a novel CR circuit.

343 CR information was redundantly encoded in the dmPFC. The advantage of the redundancy
344 is not clear, but because fear memory is critical for animal survival, it is possible that the redundant
345 coding for the fear memory is not inefficient but rather evolutionarily crucial. On the other hand,
346 although the redundancy can also be considered inefficient in terms of the short-term cost, because

347 the dmPFC is known to be involved in long-term memory via brain-wide networks (*12, 33, 34*), it

348 would be interesting to investigate whether the redundantly encoded information for the CR is

349 maintained or diminishes by longer-term continuous recording, and whether it is related to the

350 brain-wide regulation of memory using virus-based anterograde or retrograde fluorescent labeling

351 techniques to simultaneously dissect the downstream or upstream structures.

352     As we have successfully discriminated the specific neural population encoding the CR as

353 well as the detailed internal structure with a hub of the US-responsive neurons, further testing the

354 causality of the identified structure to behavior by holographic optogenetics (*30*) could be intriguing.

355 But importantly, we also found that the dmPFC responds to auditory signals even before memory

356 consolidation (Figs. 2D-F, S1) and that the CR ensemble predominantly includes the US-responsive

357 neurons (Fig. 4F). Because enhancing the sensory coding can modify behavioral responses in a task

358 based on the sensory stimuli as demonstrated before (*30*), and because activating US-responsive

359 neurons may sufficiently encourage defensive freezing behaviors as unconditioned responses,

360 further mathematical dissection and additional anatomical dissection discussed in the preceding

361 paragraph would be the next important step to more precisely identify the memory-specific

362 connections and information flow to be tested by the holographic optogenetics.

363

381

382   **Author contributions:** M.A. conceived and coordinated the whole project. M.A. designed and

383   performed behavioral experiments with the support of A.K., H.H., and T.N.; M.A. constructed in

384   vivo imaging system and performed imaging experiments with the support of Y.A. and T.N.;

385   M.A. performed data analyses with the support of I.S, Y.R.T., L.C.-R., M.Y., T.I., and J.N.; M.A.

386   wrote the paper, with contributions from all authors.

387

388   **Competing interests:** Authors declare that they have no competing interests.

389

390   **Data and materials availability:** The data that support the findings of this study are available

391   from the corresponding author upon reasonable request. Custom codes used to analyze data in this

392   study are available from the corresponding author upon reasonable request.

393

394   **Materials and Methods**
395   Animals.

396        All animal experiments were carried out in accordance with the Institutional Guidance on

397   Animal Experimentation and with permission from the Animal Experiment Committee of Osaka

398   University (authorization number: 3348), or in accordance with National Institutes of Health

399   guidelines and approved by the National Institute for Physiological Sciences Animal Care and Use

400   Committee (approval number 18A102). Male C57BL/6 or PV-Cre mice (Jax: 008069) mice housed

401   under a 12-h light/dark cycle with free access to food and water were used for all experiments.

402   Behavioral experiments were performed during the dark cycle (i.e., when mice were normally

403   awake) using single-housed mice. Mice at 4–6 months of age were used for the behavioral and

404   imaging experiments.

405

406   Virus injection

407        To express GCaMP6f, a genetically encoded calcium indicator to monitor the neural activity,

408   we used a gene expression system based on the AAV vector. Viruses were injected into mice at

409   postnatal day (P) 50-120 for in vivo experiments, at least 1 month before the microprism

410   implantation, which was followed by the in vivo experiments 1–3 months after the implantation.

411   Injection procedures were performed as described previously (*29*), with some modifications.

412   During surgery, the mice were anesthetized with isoflurane (initially 2% [partial pressure in air]

413   and then reduced to 1%). A small circle (~1 mm in diameter) of the skull was thinned over the left

414   mPFC using a dental drill to mark the site for a small craniotomy. AAV1/CamKII.GCaMP6f was

415 obtained from the University of Pennsylvania Vector Core, and injected into the left mPFC (slightly

416 away from the imaging target area to avoid damaging the field of view) at three sites (depth 1.0,

417 1.5, and 2.0 mm from the pial surface, volume 375 nl/site) to cover the dorsal mPFC, over a 5-min

418 period at each depth using a UMP3 microsyringe pump (World Precision Instruments). The X-Y

419 coordinates for the injection site was usually 0.5 mm lateral to the midline and 2.0 mm rostral to

420 bregma, but if large blood vessels obstructed the position, we shifted the insertion site slightly to

421 avoid the vessels. The beveled side of the injection needle was faced to the midline so that the

422 needle could be smoothly inserted and the virus would cover the surface layers of the mPFC. We

423 designed our injection protocol (especially the volume and depth) carefully to widely cover the

424 mPFC areas, while the anatomical coordinates of the field of view for the two-photon imaging were

425 precisely targeted using the position of the pial surface and the sinus, which were usually visible

426 through the imaging window prepared as shown below, as a guide (the field of view ranged from a

427 depth of ~0.9-1.9 mm and centered at a depth of ~1.1-1.5 mm from the pial surface and the sinus).

428

429 In vivo two-photon imaging

430     In vivo two-photon imaging was performed as described previously (*19, 29*), with

431 modifications to pair with our new experimental system. At 1−3 months after the virus injection,

432 the mice were anesthetized with isoflurane (initially 2% [partial pressure in air] and reduced to 1%).

433 A titanium head plate described in a previous paper by Goldy et al. (*35*) was selected for the present

434 study to minimize the area laying over the ear and to minimize the blockage of auditory input

435 through the ear. The head plate was attached to the skull with dental cement. For the subsequent

436 microprism implantation, a square cranial window (~2.3 x 2.3 mm) was carefully made with

437 minimal bleeding above the right mPFC, the hemisphere opposite to the virus injection site. An

438 implantable microprism assembly(*19*), comprising a 2-mm right angle glass microprism (TS N-

439 BK7, 2mm AL+MgF2, Edmund) bonded to a 2x2 mm square cover glass (No.1; Matsunami) for

440 the middle position and a 4x4 or 3x4 mm glass window at the surface position of the imaging

441 window, was prepared and inserted into the subdural space within the fissure along the midline as

442 described previously(*19*) to avoid harming any nerves surrounding the mPFC network in both

443 hemispheres, allowing for visualization of the left mPFC, which was previously injected with the

444 GCaMP6f virus, through the imaging window. The area directly beneath the microprism was

445 compressed but remained intact. This insertion procedure sometimes caused a small amount of

446 bleeding that covered the imaging site, but even in that case, the imaging window became clear

447 after waiting at least a month before performing the experiments. As reported before (*19*), the mice

448  recovered quickly and displayed no gross impairments or behavioral differences compared with
449  non-implanted mice, enabling chronic imaging of the dmPFC in behaving mice.

450      The activity of dorsal mPFC neurons was recorded by imaging fluorescence changes with
451  a FVMPE-RS two-photon microscope (Olympus) and a Mai Tai DeepSee Ti:sapphire laser
452  (Spectra-Physics) at 920 nm, through a 4x dry objective, 0.28 N.A. (Olympus) or a 16x water
453  immersion objective, 0.80 N.A. (Nikon). Mean (±SE) frame rate was $8.96 \pm 0.87$ (frames/s).
454  GCaMP6f signals were detected via the band-pass emission filter (495-540nm). As the GCaMP6f
455  was expressed under the regulation of the CaMKII promoter (*21, 22*), all of the recording targets
456  were assumed to be excitatory neurons (*20*). Scanning and image acquisition were controlled by
457  FV30S-SW image acquisition and processing software (Olympus). To smoothly set the mice below
458  the objective lens for the imaging, light and minimal-duration isoflurane (2.0% for less than 2-3
459  min) anesthesia was used, and behavioral and imaging experiments were started 5 min after the
460  mice awoke and began locomoting on the running disk, which was visually confirmed via the video
461  camera (VLG-02, Baumer) under infrared light-emitting diode illumination (850nm: LDL-
462  130X15IR2-850, CCS Inc.). To detect neural activity from the same set of neurons in each mouse
463  over multiple days, the depth from the surface of the brain (dmPFC area) and configuration of blood
464  vessels and basal GCaMP6f signals in each field of view were recorded and referenced as described
465  previously (*36*).

466

467  Fear conditioning, memory retrieval, and extinction under the microscope

468      The experiments were designed according to previous studies, with some modification to
469  optimize conditions for the two-photon microscope system (*13, 14, 18*). The heads of the mice were
470  fixed under the objective lens for two-photon imaging, allowing them to run freely on the running
471  disk placed below them, and locomotion and the freezing response were measured by the rotation
472  of the running disk, as previously described (*37*).  Experiments were performed in a completely
473  dark environment to protect the detector (photo multiplier tube) for the two-photon imaging from
474  the room light. We prepared two different types of running disks to establish two different contexts,
475  as used in conventional fear conditioning experiments for head-unfixed mice (*13, 14, 18*). Disk A
476  was made of light-colored plastic with ridges from the center to the rim that the mice could grip to
477  allow them to easily rotate (and walk on) the disk (*37*). Disk A was used for adaptation (D1 and
478  D2) and for retrieval and extinction (D4). Disk B was built for the fear conditioning (D3), and
479  comprised a grid made of stainless steel bars (Fig. 1A), which was attached to a foot shock generator
480  (SGA-2010, O'HARA & CO., LTD) via an electrical slip ring so that electrical current to this

481　running disk for the foot shock (US) could be stably delivered to the mouse irrespective of whether

482　the running disk was rotating. The behavioral sessions on each day began only after the mouse was

483　constantly locomoting for more than 5 min. The running disks and the surrounding area (inside the

484　cage for the microscope) were cleaned with 70% ethanol before and after each experiment. To score

485　freezing behavior, the speed of the mouse locomotion was measured by the rotation speed of the

486　running disk (*37*), and mice were considered to be stationary (during no CS presentation) or freezing

487　(during CS+/retrieval) if no movement was detected for at least 1 s. On D1 and D2, the mice

488　underwent an adaptation session with disk A for an hour each day, to familiarize them with the

489　novel environment. On D3, the mice underwent a habituation session in context B, in which they

490　received four presentations of the CS− and CS+ alternately (total CS duration, 30 s for each trial;

491　consisting of 50-ms pips at 1 Hz repeated 30 times; pip frequency, 7.5 kHz or white-noise,

492　respectively, 80-dB sound pressure level (60-dB basal room noise produced by the air conditioning

493　system, and 20-dB for the CS)). The habituation session was immediately followed by

494　discriminative fear conditioning (*13, 14, 18*) on the same day by pairing the CS+ with a US (1-s

495　foot shock, 7 CS+–US pairings).The intensity of the foot shock was usually 0.05~0.1 mA, but when

496　mice showed no responses at all, which was probably caused by that a part of the running disk

497　became dirty or wet by mice and the foot shock might be suppressed by this during the experiment,

498　an intensity of 0.25~0.45 mA was used. The onset of the US coincided with the onset of the last

499　sound pip of each 30-s CS trial. The CS−and the CS+ trials were performed alternately (inter-trial

500　intervals, 50–150 s). On D4, conditioned mice underwent a retrieval session followed by an

501　extinction session on disk A during which they received 4 presentations of the CS− and 12

502　presentations of the CS+. During the experiment (D1-4), the mouse was continuously encouraged

503　to locomote by administering a 4-ul drop of saccharin water per 100 cm of locomoting, provided

504　through a spout placed near their mouth (*36*) so that the freezing response could be discriminably

505　detected as decreased locomotion (Fig. 1). The mice were not water-deprived. The locomotion

506　speed and timings of the tones and the foot shock were synchronously recorded with image

507　acquisition (GCaMP6f imaging in dmPFC) using NI software (Labview; National Instruments) and

508　NI-DAQ (National Instruments). The results shown in Fig. 1 show that this protocol led to the mice

509　successfully learning the CS+-US association, and show a reduction in locomotion in response to

510　the CS+, but not the CS−, and not before but only after the fear conditioning session, enabling us

511　to observe changes in neural representations in the dmPFC as a result of the fear conditioning.

512

513　<u>Imaging data analyses and statistics</u>

514  The raw images of the GCaMP6f signals in the dmPFC were processed to correct for brain
515  motion artifacts using the enhanced correlation coefficient image alignment algorithm (*38*). To
516  apply the same regions of interest (ROIs) for analyzing the images obtained across multiple days,
517  the movies from the same mouse were precisely aligned with each other using the same enhanced
518  correlation coefficient algorithm as above, while, for a local shift (shift of a few pixels in a small
519  number of neurons among all recorded cells), the corresponding ROIs were manually adjusted.

520  The ROIs for the detection of neural activity were automatically selected using a constrained
521  nonnegative matrix factorization algorithm in MATLAB as described previously (*39*), with some
522  manual adjustment. Further steps to process the GCaMP6f signals for measurements of the signal
523  change ($\Delta F/F$) of each neuron were performed as described previously (*29, 40*); although the same
524  constrained nonnegative matrix factorization package for ROI detection also provides an option for
525  signal processing that was not sufficiently optimized to analyze our data, which were obtained over
526  several days with more than 30,000 frames each day. Fluctuations in the background fluorescence,
527  which contains synchronous fractions across nearby neurons (*39, 40*), was subtracted before
528  calculating the $\Delta F/F$ of GCaMP6f signals as described previously (*29*). Briefly, a ring-shaped
529  "background ROI" was created for each ROI 2–5 pixels away from the border of each neuronal
530  ROI to a width of 30–35 pixels, and the size was adjusted to contain at least 20 pixels in each
531  background ROI after completing the following steps. From the background ROI, we removed the
532  pixels that belonged to any neuronal ROIs, and the ROIs that contained artificially added pixels
533  (black pixels added at the edge of the image due to the motion correction procedure) at any time-
534  point. We then removed the pixels that, at some time-point(s), showed signals exceeding that of the
535  neuronal ROI by two standard deviations of the difference between each background ROI pixel
536  time series and the neuronal ROI time series. The resulting background ROI signals were averaged
537  at each time-point, and a moving average of the time series was calculated. Using the moving
538  average instead of the raw background ROI signal was helpful to minimize the production of an
539  artificially large increase or decrease at each time-point due to the subtraction, which could have
540  altered the analyses of the timing of neural activations. Pixels within each neuronal ROI were also
541  averaged to give a single time course, and then the background ROI signal was subtracted. Then,
542  the $\Delta F/F$ of GCaMP6f signals of all neurons in each circuit was calculated. For most of the analyses
543  and comparisons of the results from multiple mice, the $\Delta F/F$ data were further z-normalized within
544  each experiment (same mouse, same day) as described previously (*13, 18*). On the other hand,
545  particularly for the CRF modeling used to evaluate the functional network connectivity, the spike
546  probabilities were inferred from the $\Delta F/F$ as an alternative estimate of neuronal activation using a

547 constrained sparse nonnegative calcium deconvolution method (*39*). We used the code

548 "constrained_foopsi.m" (*39*), and the parameters used in the calculation were not manually selected

549 but estimated from the data by the code. After inference of the spike probability and further

550 thresholding by two standard deviations, the obtained binominal data were further binned (bin size:

551 1 s). Importantly, the results obtained by CRF modeling were consistent with the results of the

552 coactivity analyses based on the ΔF/F (and z-normalized ΔF/F) (Fig. S7), providing substantial

553 support that the analyses based on both estimates complemented each other for the data analyzed

554 in the present study. While neurons for the analyses were initially automatically detected, neurons

555 responding to noisy signals with no apparent calcium transient at any time during the experimental

556 days were identified by visual inspection and excluded from further analysis.

557       For the statistical analysis, we used MATLAB (MathWorks, Natick, MA). The Wilcoxon

558 signed rank tests for paired comparisons or the Wilcoxon rank sum test (equivalent to Mann-

559 Whitney U test) for unpaired comparisons was used to determine statistical significance ($P < 0.05$)

560 unless otherwise indicated. Two-tailed tests were selected for all statistical analyses. All p-values

561 less than 0.0001 are described as "$P<0.0001$" (or ****). Graphs were produced by MATLAB

562 (MathWorks) or Excel (Microsoft). When comparing two groups (e.g. D3 vs D4) consisting of the

563 results of multiple mice, in addition to the analyses using original data (e.g. N=7 vs N=7 [D3 vs

564 D4]), we performed bootstrap resampling to more systematically estimate representative values

565 (e.g. mean or median) of each mouse or each group where the number of recorded neurons in each

566 field view varied. When statistically comparing original data (e.g. comparing D3 vs D4), we used

567 a paired permutation test that does not require any assumptions regarding the data distribution,

568 though the p-values obtained by this method and the evaluated statistical significance were very

569 similar to those obtained by the paired t-test in almost all cases. For the analyses based on bootstrap

570 resampling followed by statistical comparison, random resampling (with accepting overlapped

571 sampling) from each mouse was performed in total with the same number as that of the original

572 data of each mouse for each resampling round, and the means (e.g. of 7 mice each day) and the

573 means of the difference or ratio (e.g. difference between D3 vs D4 averaged over mice) were

574 calculated. This was repeated 2000 times to derive the distribution (of 2000 bootstrap replications)

575 for each estimate, and the statistical significance was evaluated based on the 95% confidence

576 interval.

577       In the present study, to compare changes in neural responses and ensemble representations

578 before and after the fear memory consolidation without any bias, we did not exclude neurons that

579 showed no response to the CS on D4 from the analyses, which was done in some previous

580    experiments (e.g. (*18*)). Neurons for the analyses were automatically selected based on the neural

581    responses, as described above, and all neurons that exhibited clear activity during at least one of

582    the experimental days were included for the analyses irrespective of whether it was during the CS

583    presentation or only during no CS presentation, considering the previous work suggesting that not

584    only the neurons that typically respond to the CS, but also other types of neurons (including those

585    of mixed selectivity)  are important for population coding in the prefrontal network (*23*).

586         The significance of CS-induced neural responses was determined according to previous

587    studies (*13, 18*). Signals during CS presentation were normalized to baseline activity using a z-

588    score transformation, as described previously (*13, 18*). The CS-induced neural activity for each

589    stimulus was then calculated as the mean of the activity during ~1 s from each stimulus onset

590    (depending on the imaging frame rates, we set the number of frames to be used for this calculation

591    so that sampling duration was closer to 1 s but the frames that overlapped with the next stimulus

592    onset was excluded). The last sound pip of each 30-s CS trial was also excluded from this analysis

593    because, during fear conditioning, the last sound pip of the CS+ overlapped with the US (we

594    excluded the last pip data not only for analysis of CS+-evoked responses during fear conditioning

595    but for all data analyses on both D3 and D4, for both CS+ and CS−). They were averaged over

596    blocks of 3 CS trials consisting of 87 individual sound pips in total, for D3E (first three trials during

597    the fear conditioning session), D3L (last three trials during fear conditioning on D3), D4E (first

598    three trials on D4, as responses during fear memory retrieval), and D4L (last three trials only for

599    CS+ on D4 as responses during extinction), respectively, or used to statistically test whether the

600    responses of each neuron were significantly different from zero (baseline) and to define CS-

601    activated / -inactivated neurons.

602         To define US responsive neurons, because the number of US were limited (7 stimuli in total

603    for each mouse), the mean z-score of each neuron for 1.5 s from the US onset was calculated, and

604    US responsive neurons were defined as neurons with responses of one standard deviation or larger.

605    The number of USR was very limited (zero or only a few for some of the mice) as they were only

606    around 5 % on average, and therefore all the analyses shown in Figs. 4I-K were performed with

607    pooled data from all mice (N=7 mice).

608         To evaluate the coactivation of neural activity in the dmPFC network, we calculated cell-to-

609    cell pair-wise correlations within each ensemble using Pearson's correlation coefficient, from the

610    GCaMP6f signals (z-normalized ΔF/F) of two cells over the duration of the CS+ presentation, as

611    described before (*29*). The calculated correlation coefficients (R) were statistically analyzed. As a

612    complementary analysis, we also used the inferred spike probability to analyze the functional

613 connectivity, as explained in the section describing the CRF model, which revealed consistent
614 results as shown in the results section. We further performed analyses based on surrogate datasets,
615 as described in previous studies (*29, 41*). For this, the total activity of each neuron was preserved,
616 but only the timing was shuffled randomly within each neuron, followed by calculation of the
617 correlation coefficients of shuffled data.

618

619 Extraction of neuronal ensembles

620     To directly differentiate neural populations (ensembles) encoding the CR (i.e., suppressed
621 locomotion triggered by CS+ during the memory retrieval) and those encoding RL (i.e., stationary
622 or locomotive state during no CS presentation), we used the elastic net(*28*), a regularization and
623 variable selection algorithm that enabled us to systematically extract neurons encoding respective
624 target behaviors. For this, we used the "lassoglm" function of MATLAB R2019b. Because this
625 method allowed us to identify different ensembles for different behaviors independently from the
626 same mice, we used this to verify whether neurons in CR ensembles were unique or mostly
627 overlapped with RL ensembles (Figs. 3 and S2). Compared with the conventional sparse modeling
628 method called Lasso (least absolute shrinkage and selection operator), the advantage of the elastic
629 net is that the hyper parameter "alpha" additively enables the adjustment of the size of selected
630 neurons depending on the data; when the analyzed data include strongly correlated pairs, which
631 appeared to be the case for our data as shown in Fig. S7, conventional Lasso removes redundant
632 predictors and selects only one or a part of such a synchronous population, but in the elastic net,
633 lowering the alpha value increases their inclusion, which is helpful toward preventing missing
634 encoder neurons.

635     When extracting the CR ensemble, we used data only during the CS+ presentation of D4E
636 (retrieval session) and identified neurons informative for distinguishing whether animals exhibited
637 freezing behavior or were locomoting during the CS+ so that the auditory information of the CS
638 was not considered for identifying the ensemble neurons. While mice exhibited the CR as
639 suppressed locomotion during the fear memory retrieval session (Fig. 1), they also showed more or
640 less locomotion intermittently, and both labels (freezing and locomotive) are required to perform
641 the regression based on the elastic net; only the data containing at least 10% of each label (freezing
642 and locomotive) were used to discriminate ensembles in the present study. On the other hand, for
643 extracting the RL ensemble, we used data only during the no-CS presentation (for D3 and D4).
644 Learning the elastic net is formulated as follows.

$$\min_{\beta_0, \beta} \left( \frac{1}{2N} \sum_{i=1}^{N} \left( -y_i \log \tilde{y}_i - (1 - y_i) \log(1 - \tilde{y}_i) \right) + \lambda P_\alpha(\beta) \right),$$

where

$$\tilde{y}_i = \frac{1}{1 + \exp(-(\beta_0 + x_i^\top \beta))} \quad (i = 1, \ldots, N)$$

$$P_\alpha(\beta) = \frac{(1 - \alpha)}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 = \sum_{j=1}^{p} \left( \frac{(1 - \alpha)}{2} \beta_j^2 + \alpha \,|\beta_j| \right)$$

and N is the number of observations; $y_i$ is the behavior (freezing/stationary y_i=1 or locomotive y_i =0) at observation i; $x_i$ is data (neuronal activity), a vector of p values at observation i; $\lambda$ is a positive regularization parameter; parameters $\beta_0$ and $\beta$ are a scalar variable and a p-dimensional vector, respectively. As $\lambda$ increases, the number of nonzero components of $\beta$ decreases. The elastic net is a hybrid of ridge regression and lasso regularization: when alpha ($\alpha$) = 1, elastic net is the same as lasso, while, as $\alpha$ shrinks toward 0, elastic net approaches ridge regression. For other values of alpha ($\alpha$), the penalty term P$\alpha$($\beta$) interpolates between the $L^1$ norm of $\beta$ and the squared $L^2$ norm of $\beta$. Lasso is sensitive to correlations between variables and can choose one if there are two highly correlated and useful variables, whereas elastic net is more likely to select both useful variables, which leads to more stable variable selection. The tuning parameter $\lambda$ controls the overall strength of the penalty. $\beta_j$ is the coefficient for the corresponding neuron j estimated by this model. Because this method is designed to sparsely leave the coefficients $\beta_j$ for the respective neurons, we could identify neurons with a non-zero coefficient as ones of substantial decodability (i.e., ensemble neurons). The lambda value with minimum expected deviance, as calculated during cross-validation, was selectively used to define these beta coefficients for each dataset. To avoid an imbalance of the number of original labels for respective states (e.g. freezing or locomotive for CR ensembles) for the training, the same number of data points from respective states were randomly selected to prepare the training data despite an overlap, a total of 900 samples for each, and used to produce the model. We found that the eventual model and non-zero-coefficient neurons slightly varied trial by trial. To accurately define each ensemble, we repeatedly performed this procedure (random sampling and modeling) 100 times to obtain the distribution of each beta value. Gaussian fitting was performed to define the centroid and the 95% confidence interval of each distribution of each beta, and then the 95% confidence interval was used to determine whether or not they were significantly different from zero (enabling us to maintain sparsity), with the centroid being used to define the final beta values of non-zero coefficient neurons to build the model. To evaluate the fitting and decoding performance of the obtained model, the prediction accuracy and the area under

674 the curve (AUC) of receiver operating characteristic curve (ROC) were calculated, respectively,

675 revealing that those scores were very similar and highly correlated with each other (Fig. S5).

676       Based on the above-described procedure, we next optimized the alpha values. Ideally, if all

677 the informative neurons can be extracted into the selected CR ensembles, the remaining neurons

678 should have poor decoding performance. According to this idea, to optimize the alpha value, after

679 building a model at each alpha for each mouse ("AUC original" in Fig. S3A), we compared the

680 difference in decoding performance between "AUC CRE-rem" and "AUC nonCRE-rem" (Figs.

681 S3A). AUC CRE-rem is the AUC value calculated by an elastic net model built with the neurons,

682 excluding the original CR ensemble neurons. On the other hand, AUC nonCRE-rem is the AUC

683 value calculated by the neurons, excluding neurons other than original CR ensemble neurons,

684 randomly selected, and the number of excluded neurons was the same as the number of original CR

685 ensemble neurons (so that the number of neurons used to calculate AUC nonCRE-rem were set to

686 be the same as that used for AUC CRE-rem calculation). The "AUC difference" (Fig. S3A) between

687 those two values was calculated to estimate the degree of remaining information, and in principle,

688 we defined the best alpha based on the maximum AUC difference for each mouse independently.

689 In addition, for further statistical evaluation to define the optimal alpha as explained below, we

690 repeated these procedures 10 times for both "AUC CRE-rem" and "AUC nonCRE-rem".

691       As shown in Fig. S3B, although the decoding performance of the original CR ensembles

692 (i.e., AUC original in Fig. S3A) was not affected by the alpha (Fig. S3B, middle), the size of the

693 CR ensemble was affected, and a smaller alpha generally resulted in a larger number of selected

694 neurons for each CR ensemble (Fig. S3B, top), suggesting that the CR information might be

695 redundantly encoded in the dmPFC as discussed in detail later. On the other hand, the influence of

696 the alpha on the AUC difference was more complicated. As explained above, we defined the best

697 alpha based on the maximum AUC difference for each mouse independently, but in some

698 exceptional cases as shown in Fig. S3D (mouse #3), when the other alpha(s) showed a AUC

699 difference(s) not significantly far from the maximum AUC difference, the alpha of the smallest of

700 the ensembles among those alphas, i.e., largest alpha among them, was selected to avoid

701 unnecessarily including additional neurons that did not improve the AUC difference (e.g. in mouse

702 #3, alpha = 0.1, 0.05, 0.01 showed similar AUC differences and there was no statistically significant

703 difference between them [Wilcoxon rank sum test, alpha of maximum AUC difference vs the other

704 alpha, n=10 estimates for each calculated as explained above], so in this case, the largest alpha 0.1

705 among those three was selected to define the CR ensemble for this mouse).

706      These results revealed two important points. First, searching around the alpha value may be

707    important in some cases. Considering this, we also searched alphas in the case of RL ensembles

708    (Fig. S4), and found that there was no difference among the various alphas, for the RL ensembles,

709    even if we tested an additional number of reference frames (means of the neural activities over the

710    past or future several frames were used as neural activity data to predict a single label at each single

711    time-point, which showed no significant difference from each other, evaluated by the Friedman test,

712    a non-parametric statistical test similar to the parametric one-way repeated measures ANOVA).

713    Therefore, in the present study, we fixed the alpha to define RL ensembles at 0.75 for most of the

714    analyses, except for the data in Figs. S4 and S5, where we evaluated the influence of the alpha for

715    RL ensembles.

716      Second, fear memory triggering the CR might be redundantly encoded in the dmPFC. As

717    discussed above, although decoding performance of the original CR ensembles was not affected by

718    the alpha (Fig. S3B, middle), the size of the CR ensemble was affected, and a smaller alpha

719    generally resulted in a larger number of selected neurons for each CR ensemble (Fig. S3B, top). In

720    addition, when the alpha was fixed at alpha (A) =0.9 (a larger alpha (than 0.9) did not work for

721    some circuits in our data), while the uniqueness of the CR ensembles was maintained and the ratio

722    of the CR ensemble neurons overlapping with RL ensembles was 26.84% (Fig. S3E), which was

723    very similar to the case of alpha-optimized CR ensembles (Fig. 3), the size of this CR ensemble

724    (A=0.9) was two times smaller than that of the alpha-optimized CR ensembles (Fig. S3F).

725    Importantly, 97.82% of the neurons selected at A=0.9 were also selected in the alpha-optimized CR

726    ensembles (Fig. S3F), suggesting that the neurons selected at the largest alpha might be more

727    reliable and robust for the decoding among all the informative neurons. In addition, even after the

728    removal of such "core" neurons, the remaining neurons also possessed information for the CR (Figs.

729    S3B, D), indicating that the CR information was redundantly encoded in the dmPFC. Because this

730    redundancy was specific to the CR ensemble and not observed in the RL ensemble, it would be

731    interesting to investigate possible changes in this redundancy when the memory is recalled as a

732    long-term memory (e.g. 30 days after the memory consolidation).

733      To evaluate the dominance of the CR ensembles vs the RL ensembles, we applied the CR

734    decoder to predict the RL, and vice versa (Figs. 3 and S5).

735

736    <u>CRF models to evaluate functional connectivity</u>

737      To evaluate the functional connectivity between neurons in the recorded network and the

738    pattern completion capability of each neuron, we used conditional random fields (CRFs) as

739 described previously (*30*), which models the conditional probability distribution of a given neuronal

740 ensemble firing together. We used CRFs to capture the contribution of specific neurons to the

741 overall network activity defined by population vectors belonging to a given neuronal ensemble. We

742 generated a graphical model in which each node represents a neuron in a given ensemble and edges

743 represent the dependencies between neurons. For training, 80% of the recorded data randomly

744 selected from all time frames was used, and for cross-validation, the remaining 20% was used. For

745 this analysis, binned neural activity data (1 s) were used. The model parameters were determined

746 by the local maximum of the likelihood function in the parameter space. We constructed a CRF

747 model in two steps: (1) structure learning, and (2) parameter learning. For the structure learning,

748 we generated a graph structure using ℓ1-regularized neighborhood-based logistic regression(*31*).

749 Here λs is a regularization parameter that controls the sparsity (or conversely, the density) of the

750 constructed graph structure, leaving only relevant functional connectivity, including both coactive

751 and suppressive relationships. A previous study showed that this number of connections was

752 enhanced as a result of optogenetic rewiring of the local network(*31*), demonstrating the reliability

753 of the functional connectivity estimated by CRFs models. Therefore, we also calculated the ratio of

754 these remaining connections per all the possible connections for each neuron as a "functional

755 connectivity" score for each node, after carefully screening the optimal λs value by maximizing the

756 log-likelihood of the observations at the following parameter learning step. When comparing the

757 connectivity between different ensembles (e.g. within-CR-ensemble vs within-Non-CRE) or

758 different cell types (e.g. USR vs non-US responsive neurons), we first calculated a whole network

759 connectivity without separating the ensembles, and further separated them into different categories.

760 To measure which neurons were the most informative for a given stimulus (CS+ or CS−), we

761 computed the standard ROC, taking as ground truth the timing of a particular CS. The AUC from

762 the ROC curve that represents the performance of each neuron was calculated to compare the

763 encoded information in different ensembles, different neuron types, and different days (e.g. before

764 vs after the fear memory consolidation). As was recently demonstrated (*30*), high ranks for this

765 value indicate high potential to recall the neural and cognitive representation of a given stimulus.

766

767 <u>Statistical Analysis</u>

768       Statistical analyses in the present study were performed as described above (in "Materials

769 and Methods" as well as in the main text and figure legends). The data that support the findings of

770 this study are available from the corresponding author upon reasonable request. Custom codes used

771 to analyze data in this study are available from the corresponding author upon reasonable request.

772

773

## References

775

1. P. I. Pavlov, Conditioned reflexes: An investigation of the physiological activity of the cerebral cortex. *Ann Neurosci* **17**, 136-141 (2010).

2. J. B. Watson, R. Rayner, Conditioned emotional reactions. *Journal of Experimental Psychology* **3**, 1-14 (1920).

3. S. A. Josselyn, S. Tonegawa, Memory engrams: Recalling the past and imagining the future. *Science* **367**, (2020).

4. E. K. Miller, J. D. Cohen, An integrative theory of prefrontal cortex function. *Annu Rev Neurosci* **24**, 167-202 (2001).

5. G. G. Calhoon, K. M. Tye, Resolving the neural circuits of anxiety. *Nature neuroscience* **18**, 1394-1404 (2015).

6. P. Le Merre, S. Ahrlund-Richter, M. Carlen, The mouse prefrontal cortex: Unity in diversity. *Neuron* **109**, 1925-1944 (2021).

7. L. M. Shin, I. Liberzon, The Neurocircuitry of Fear , Stress , and Anxiety Disorders. *Neuropsychopharmacology* **35**, 169-191 (2009).

8. A. Jezzini, E. S. Bromberg-Martin, L. R. Trambaiolli, S. N. Haber, I. E. Monosov, A prefrontal network integrates preferences for advance information about uncertain rewards and punishments. *Neuron* **109**, 2339-2352 e2335 (2021).

9. A. Burgos-Robles *et al.*, Amygdala inputs to prefrontal cortex guide behavior amid conflicting cues of reward and punishment. *Nature neuroscience* **20**, 824-835 (2017).

10. R. J. Fenster, L. A. M. Lebois, K. J. Ressler, J. Suh, Brain circuit dysfunction in post-traumatic stress disorder: from mouse to man. *Nature reviews. Neuroscience* **19**, 535-551 (2018).

11. K. A. Corcoran, G. J. Quirk, Activity in prelimbic cortex is necessary for the expression of learned, but not innate, fears. *The Journal of neuroscience : the official journal of the Society for Neuroscience* **27**, 840-844 (2007).

12. F. H. Do-Monte, K. Quinones-Laracuente, G. J. Quirk, A temporal shift in the circuits mediating retrieval of fear memory. *Nature* **519**, 460-463 (2015).

13. J. Courtin *et al.*, Prefrontal parvalbumin interneurons shape neuronal activity to drive fear expression. *Nature* **505**, 92-96 (2014).

14. C. Dejean *et al.*, Prefrontal neuronal assemblies temporally control fear behaviour. *Nature* **535**, 420-424 (2016).

15. O. Klavir, M. Prigge, A. Sarel, R. Paz, O. Yizhar, Manipulating fear associations via optogenetic modulation of amygdala inputs to prefrontal cortex. *Nature neuroscience* **20**, 836-844 (2017).

16. D. Jercog *et al.*, Dynamical prefrontal population coding during defensive behaviours. *Nature* **595**, 690-694 (2021).

17. A. Burgos-Robles, I. Vidal-Gonzalez, G. J. Quirk, Sustained conditioned responses in prelimbic prefrontal neurons are correlated with fear expression and extinction failure. *The Journal of neuroscience : the official journal of the Society for Neuroscience* **29**, 8474-8482 (2009).

18. C. Herry *et al.*, Switching on and off fear by distinct neuronal circuits. *Nature* **454**, 600-606 (2008).

19. R. J. Low, Y. Gu, D. W. Tank, Cellular resolution optical access to brain regions in fissures: imaging medial prefrontal cortex and grid cells in entorhinal cortex. *Proceedings of the National Academy of Sciences of the United States of America* **111**, 18739-18744 (2014).

20. S.-H. Lee *et al.*, Activation of specific interneurons improves V1 feature selectivity and visual perception. *Nature* **488**, 379-383 (2012).

21. S. L. Resendez *et al.*, Visualization of cortical, subcortical and deep brain neural circuit dynamics during naturalistic mammalian behavior with head-mounted microscopes and chronically implanted lenses. *Nat Protoc* **11**, 566-597 (2016).

22. T.-W. Chen *et al.*, Ultrasensitive fluorescent proteins for imaging neuronal activity. *Nature* **499**, 295-300 (2013).

23. M. Rigotti *et al.*, The importance of mixed selectivity in complex cognitive tasks. *Nature* **497**, 585-590 (2013).

24. S. Reinert, M. Hubener, T. Bonhoeffer, P. M. Goltstein, Mouse prefrontal cortex represents learned rules for categorization. *Nature* **593**, 411-417 (2021).

25. S. Fusi, E. K. Miller, M. Rigotti, Why neurons mix: high dimensionality for higher cognition. *Current opinion in neurobiology* **37**, 66-74 (2016).

26. R. R. Rozeske *et al.*, Prefrontal-Periaqueductal Gray-Projecting Neurons Mediate Context Fear Discrimination. *Neuron* **97**, 898-910 e896 (2018).

27. K. Ghandour *et al.*, Orchestrated ensemble activities constitute a hippocampal memory engram. *Nat Commun* **10**, 2637 (2019).

28. H. Zou, T. Hastie, Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **67**, 301-320 (2005).

29. M. Agetsuma, J. P. Hamm, K. Tao, S. Fujisawa, R. Yuste, Parvalbumin-Positive Interneurons Regulate Neuronal Ensembles in Visual Cortex. *Cerebral cortex* **28**, 1831-1845 (2018).

30. L. Carrillo-Reid, S. Han, W. Yang, A. Akrouh, R. Yuste, Controlling Visually Guided Behavior by Holographic Recalling of Cortical Ensembles. *Cell* **178**, 447-457 e445 (2019).

31. L. Carrillo-Reid *et al.*, Identification of Pattern Completion Neurons in Neuronal Ensembles using Probabilistic Graphical Models. *The Journal of Neuroscience*, (2021).

32. D. O. Hebb, *The Organization of Behavior: A Neuropsychological Theory*. (Wiley, 1949).

33. T. Kitamura *et al.*, Engrams and circuits crucial for systems consolidation of a memory. *Science* **356**, 73-78 (2017).

34. P. Tovote, J. P. Fadok, A. Luthi, Neuronal circuits for fear and anxiety. *Nature reviews. Neuroscience* **16**, 317-331 (2015).

35. G. J. Goldey *et al.*, Removable cranial windows for long-term imaging in awake mice. *Nat Protoc* **9**, 2515-2538 (2014).

36. Y. Masamizu *et al.*, Two distinct layer-specific dynamics of cortical ensembles during learning of a motor task. *Nature neuroscience* **17**, 987-994 (2014).

37. S. Inagaki *et al.*, Imaging local brain activity of multiple freely moving mice sharing the same environment. *Sci Rep* **9**, 7460 (2019).

38. G. D. Evangelidis, E. Z. Psarakis, Parametric image alignment using enhanced correlation coefficient maximization. *IEEE transactions on pattern analysis and machine intelligence* **30**, 1858-1865 (2008).

39. E. A. Pnevmatikakis *et al.*, Simultaneous Denoising, Deconvolution, and Demixing of Calcium Imaging Data. *Neuron* **89**, 285-299 (2016).

40. A. J. Peters, S. X. Chen, T. Komiyama, Emergence of reproducible spatiotemporal activity during motor learning. *Nature* **510**, 263-267 (2014).

41. J.-e. K. Miller, I. Ayzenshtat, L. Carrillo-Reid, R. Yuste, Visual stimuli recruit intrinsically generated cortical ensembles. *Proceedings of the National Academy of Sciences of the United States of America* **111**, E4053-4061 (2014).

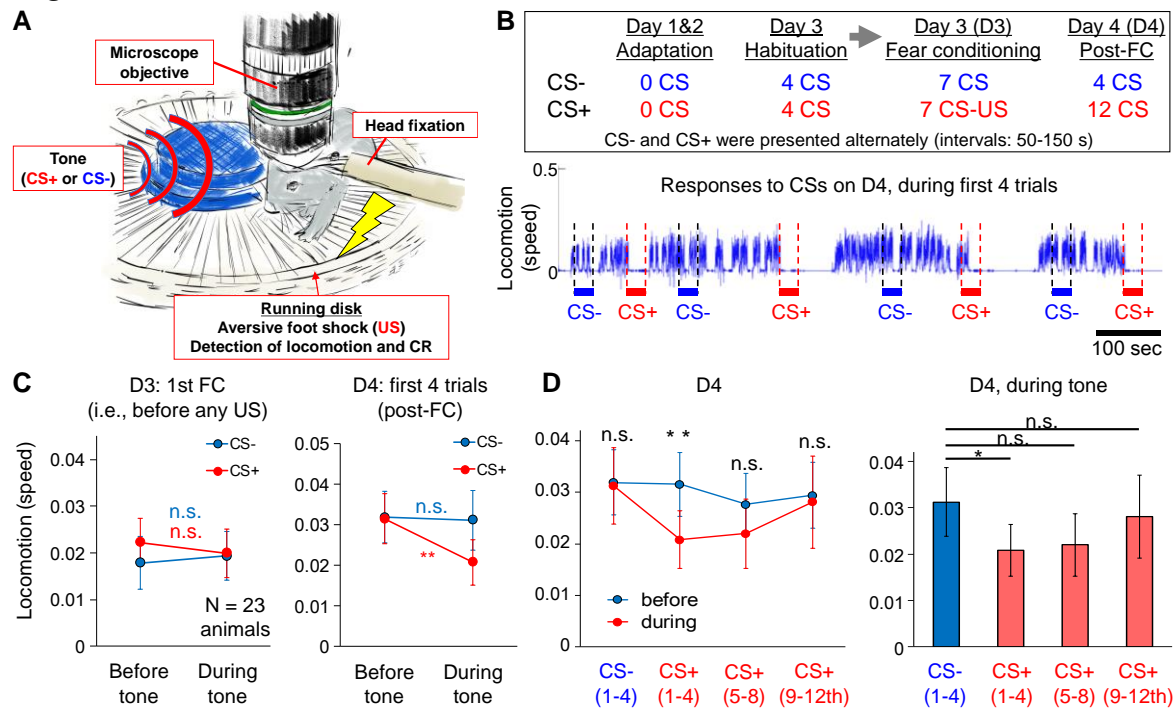872
873
874 **Figures and Tables**
875

## Fig. 1



**Fig. 1. Cued fear conditioning during two-photon microscopy.** (**A**) Schematic diagram showing the system used to perform the cued-fear conditioning and memory retrieval under a two-photon microscope. (**B**) (top) Experimental protocol. CS, conditioned stimulus; US, unconditioned stimulus; FC, fear conditioning. (bottom) An example of the changes in locomotion over time of a mouse on day (D) 4 (first four trials). (**C-D**) Fear conditioning under the microscope produced CS+-specific memory consolidation. Comparisons of the locomotor speed between before the tone onset and during the tone presentation are shown in (C-D). Before the fear conditioning (on D3), the mice (N = 23) exhibited no significant change in locomotion during the CS+ and CS- presentations (C, left, and D). After the fear conditioning (i.e., during fear retrieval; the first four trials on D4), however, the CS+ suppressed locomotion as a CR, while the CS- induced no significant change (C, right, and D). After repeated presentations of the CS+ (fear extinction; 5th-12th trials on D4), the CS+-evoked CR became smaller until no significant change in locomotion was observed upon CS+ presentation (D). (**E**) Statistical comparison among responses to the CS- and those to the CS+ at each testing phase on D4 during the tone presentation revealed that locomotion during CS+ was significantly lower only during trials 1–4 on D4, and not after repeated presentations to the CS+ (5th-12th trials). Note that locomotion during pre-tone-onset (before) was not significantly different between the CS- and CS+ conditions. *p<0.05; **p<0.01; n.s., not significant by Wilcoxon signed-rank test (the Friedman test followed by post-hoc multiple comparisons revealed similar results for panel E). Error bars, s.e.m.
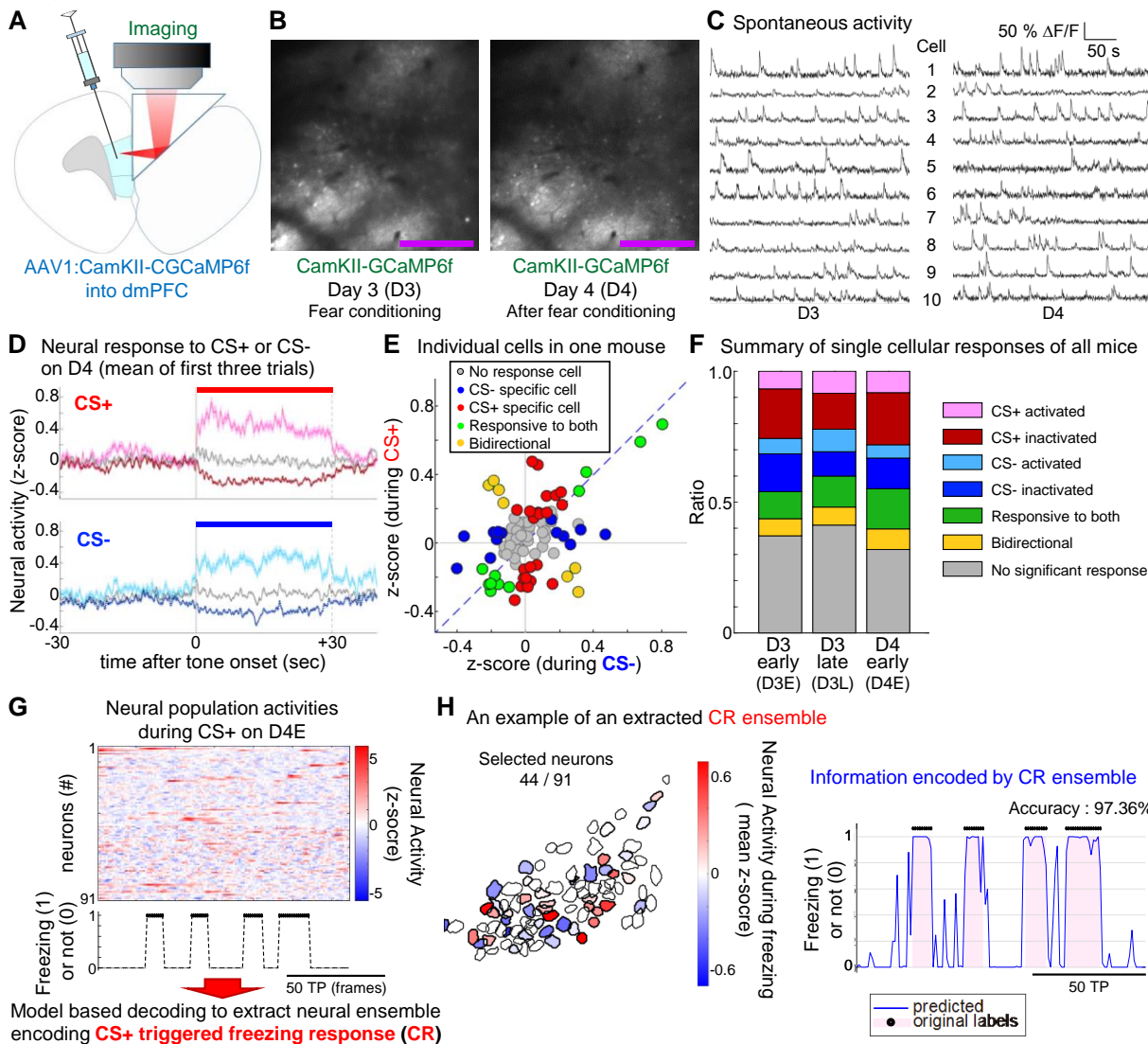
## Fig. 2



900

901 **Fig. 2. Longitudinal in vivo imaging in dmPFC and extraction of the neural ensemble**
902 **encoding conditioned responses.** (**A**) Microprism implantation along the midline for optical
903 access to the dmPFC without cutting nerves. To visualize activity of excitatory neurons in the
904 dmPFC, GCaMP6f was expressed by the AAV under the regulation of the CamKII promoter. (**B**)
905 In vivo two-photon microscopy to detect single-cell neural activity visualized by GCaMP6f,
906 chronically (day [D] 3 and D4) from the same set of neurons observed through the prism. See also
907 Movie S1 and 2. Scale bar, 250 μm. (**C**) Traces of spontaneous $Ca^{2+}$ activity from 10 example
908 neurons in dmPFC, chronically on D3 and D4. (**D**) Summary of neural responses during the retrieval
909 session (D4-early [D4E], mean of first three trials) to the CS+ or CS-. Mean of neural responses in
910 each category (significantly activated [bright red or blue], inactivated [dark red or blue], and others
911 [dark gray]), as well as the mean of all cells (light gray) are plotted. (**E**) Scatter plot showing
912 responses of individual neurons to the CS+ and CS- in an example mouse during D4E. Each dot
913 represents the mean response of each neuron. Blue, red, and green colors indicate that cells had a
914 significant response as described in the panel. These features for all the mice are summarized in
915 panel E. (**F**) Summary of response profiles at each phase (D3E, D3-late [D3L], and D4E,
916 respectively; N=7 chronically recorded mice). (**G**) Schematic diagram showing how we extracted
917 the CR ensembles. See the Materials and Methods for details. (**H**) An example of the CR ensemble

918    and encoded neural representation of the behavior. (left) Extracted neurons are drawn with a bold
919    margin, and the mean activity during CR (freezing) is shown in color. (right) Time course changes
920    of neural representation encoded by the CR ensemble shown in the left panel. Black dots on the top
921    of the graph and pink color in the graph indicate the timing of the actual CR, while the blue line
922    shows information decoded by this CR ensemble. The plots show a part of the whole length of the
923    data, and overall decoding accuracy was 97.36% in this example. TP, time points (i.e., image
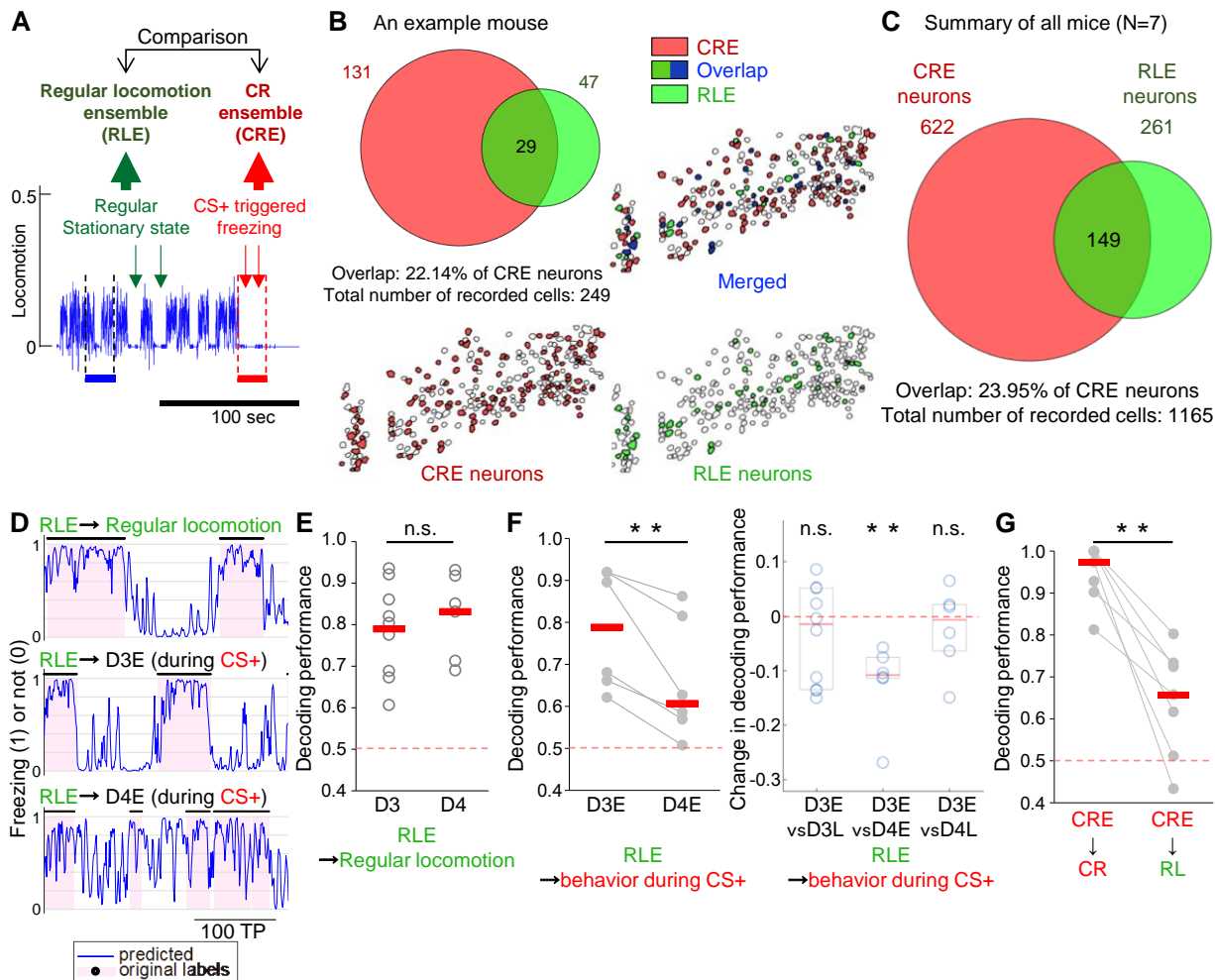924    frames).

925

## Fig. 3



926
**Fig. 3. Emergence of unique CR ensembles after fear conditioning.** (**A**) The CR ensemble was
compared with the RL ensemble to evaluate the uniqueness. (**B**) An example Venn diagram and an
example spatial map showing the overlap between CR ensemble neurons and RL ensemble neurons
in an example mouse. (**C**) Summary of the overlap between the CR ensemble neurons and RL
ensemble neurons of all mice (N=7, n=1165 neurons). (**D-F**) Decoding locomotion during regular
locomotion (inter-trial interval) or during CS+ by the RL ensemble. (**D**) In an example mouse, an
RL ensemble (RLE) that showed high accuracy for decoding performance to predict RL (top) also
showed high decoding performance in predicting locomotion during CS+ at day 3-early (D3E). But
the performance dropped when locomotion during CS+ at D4E (i.e., during fear retrieval) was also
predicted by the RL ensemble. (**E**) Original decoding performance of the RL ensembles (i.e.,
predictability for RL) were not significantly different between D3 and D4. (**F**) (left) Decoding
performance of RL ensembles to locomotion during CS+ at D4E (i.e., during fear retrieval) was
significantly lower than that for D3E (i.e., before memory consolidation). (right) The change in
decoding performance was systematically evaluated. Decoding performance was not significantly
different between D3E and D3-late (D3L), or between D3E and D4L. (**G**) Decoding locomotion
during CS+ by CR ensembles. Decoding performance was significantly decreased when the CR
ensembles were applied to predict RL. Within D3, N=10; D3 vs D4 and within D4, N=7 pairs. A
non-paired comparison (Wilcoxon rank sum test) was performed for panel D, while for the other
comparisons in E and F, a paired permutation test was performed. For the decoding performance,
we plotted the accuracy scores, while the AUC was very similar as shown in Fig. S5. **p<0.01;
n.s., not significant. Red bars, median; box in panel E (left) indicates 25th and 75th percentiles.
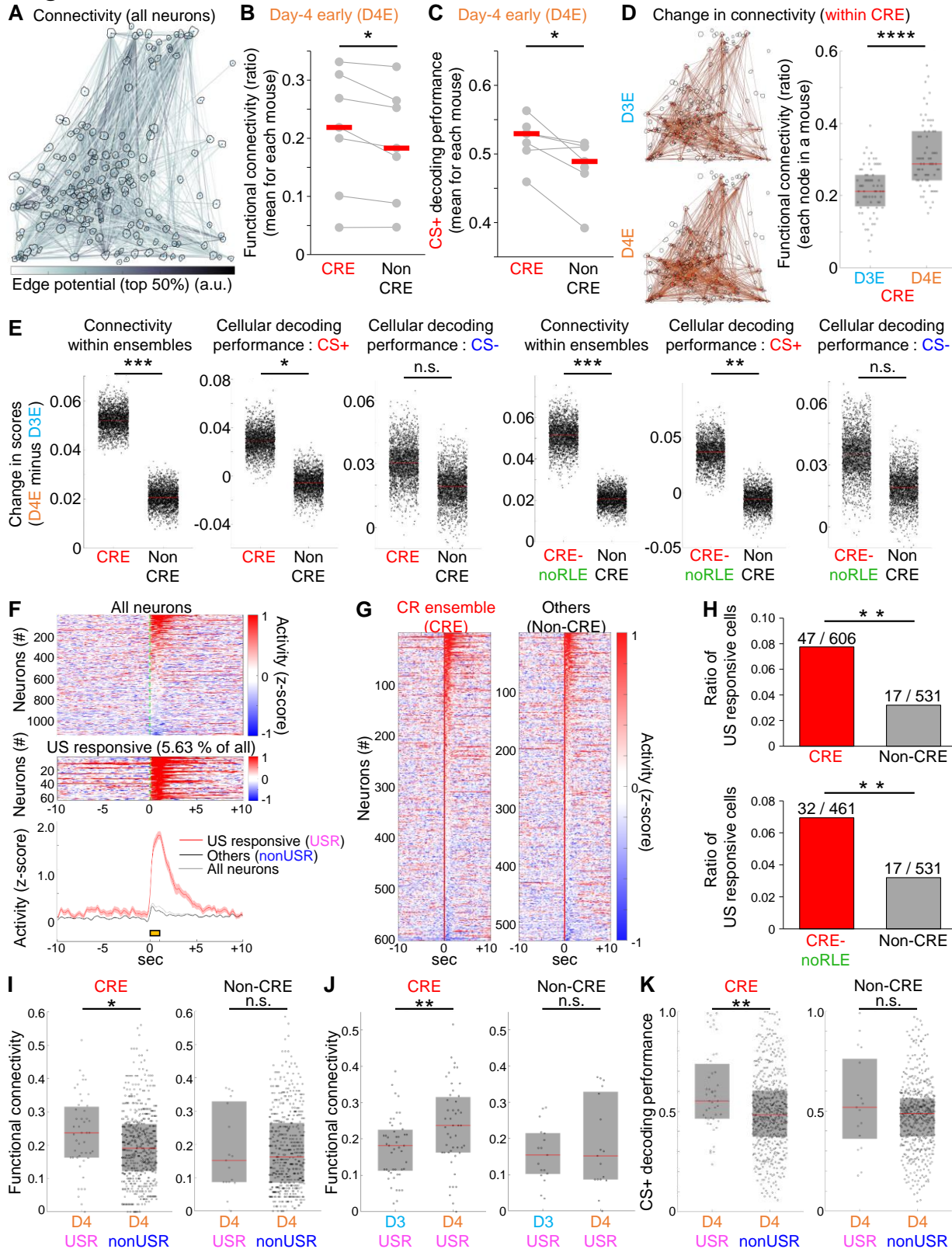
# Fig. 4

**Fig. 4. Enhanced functional connectivity and CS+ predictability in the CR ensemble with an emergent hub of US responsive neurons after fear conditioning.** (**A**) Functional connectivity between neurons in an example circuit. Among all the possible connections for all pairs of neurons, the CRF model enables the estimation of functional connections, as well as the dependencies of connected pairs. In this panel, the top 50% edge potentials were visualized. (**B**) During day 4-early [D4E], the functional connectivity within CRE was significantly higher than that of Non-CRE. (**C**) During D4E, the predictability for CS+ in CRE was also significantly higher than that in Non-CRE. (**D**) Change in functional connectivity within CRE of an example circuit. This is the same as the circuit shown in A, but only the connectivity of the CRE neurons marked by the red ellipses were analyzed. Left panel shows the change in the connectivity between day 3-early [D3E] and D4E, while the right panel shows the change in the ratio of functional connectivity per all possible connections for individual nodes (i.e., individual neurons). (**E**) Summary of changes in functional connectivity and cellular decoding performance for CS+ and CS- of all observed networks (N=7 mice). Differences (D4E minus D3E) of these scores are plotted as a result of bootstrap resampling (2000 times) to compare CRE and Non-CRE, or CRE-noRLE (CRE neurons excluding those overlapping with RL ensemble neurons) and Non-CRE. (**F**) A part of the recorded neurons in the dmPFC showed increased activity upon US presentation on day 3 (D3) during fear conditioning. Mean activity over 7 trials of all (top) or US-responsive (middle) neurons, and the mean ± s.e.m. of respective categories (bottom) are plotted. Green dotted line indicates the onset of the US, and yellow bar indicates the 1-s duration of the US presentation. (**G**) Summary of US responses of CR ensemble neurons (CRE) and others (Non-CRE). All individual neurons for the respective categories are plotted. (**H**) Neurons responding to the US on D3 were predominantly involved in the CRE on D4 after the fear conditioning. The difference between CRE vs Non-CRE, as well as CRE-noRLE vs Non-CRE, was statistically evaluated. (**I**) Comparison of functional connectivity between US responsive neurons (USR) and others (nonUSR) on D4. In the CRE network, USR became more connected within the network than nonUSR, while there was no significant difference between USR and nonUSR outside of the CRE (Non-CRE). (**J**) The higher connectivity of USR on D4 was experience-dependent. Functional connectivity of USR on D4 was significantly higher in CRE, while there was no significant difference between them in Non-CRE. (**K**) USR in the CRE exhibited significantly higher decoding performance of CS+ than nonUSR, which was not the case in Non-CRE. A paired permutation test was used for the statistics in B and C. The Wilcoxon signed-rank test was used for the statistics in D. The data obtained by bootstrap resampling were statistically analyzed as described in the Materials and Methods. Because the number of USR was limited (only 5.63% under the present definition), the analyses shown in panels F-K were performed with data pooled together from all mice (N=7 mice). Fisher's exact test was used for the statistics in H, a non-paired comparison (Wilcoxon rank sum test) was used in I and K, and the Wilcoxon signed-rank test was used in J.  *p<0.05; **p<0.01; ***p<0.001; ****p<0.0001; n.s., not significant. Red bars, median; gray boxes in panels D, E, I-K indicate 25th and 75th percentiles.

993

**Supplementary Materials**

995

996    Supplementary Materials (Figs. S1 to S8, and caption for Movies S1 and S2) are explained
997    in a separate document.

998