

1302 A Supplement to methods

1303 A.1 HCP experimental task details

1304 **Working Memory** Each of the two runs of the working memory task consisted of eight task
1305 (25 s each) and four fixation blocks (15 s each). In each of the task blocks, participants saw images
1306 of one of four different stimulus types (namely, images of body parts, faces, places or tools). These
1307 four stimulus types are known to reliably engage distinct cortical regions (Downing et al., 2001)
1308 across subjects Peelen and Downing (2005) and time (Fox et al., 2009). Half of the task blocks
1309 used a 2-back working memory task (participants were asked to respond “target” when the current
1310 stimulus was the same as the stimulus 2 back) and the other half a 0-back working memory task (a
1311 target stimulus was presented at the beginning of each block and participants were asked to respond
1312 “target” whenever the target stimulus was presented in the block). Each task block consisted of 10
1313 trials (2.5 s each). In each trial, a stimulus was presented for 2 s followed by a 500 ms interstimulus
1314 interval (ISI). We were not interested in identifying any effect of the N-back task condition on the
1315 evoked brain activity and therefore pooled the data of both N-back conditions.

1316 **Gambling** Participants played a card guessing game in which they were asked to guess the
1317 number on a mystery card. The potential card numbers ranged from 1 to 9 and participants
1318 were asked to indicate whether they think that the number is going to be above or below 5.
1319 Participants received feedback in form of a number on the card. Importantly, the number on the
1320 card was dependent on whether the respective trial was marked as a reward, loss, or neutral trial.
1321 In addition to the number, the feedback included a green arrow pointing upwards with “1” for
1322 reward trials or a red arrow pointing downwards next to “-0.50” for loss trials or the number
1323 “5” and a gray double headed arrow for neutral trials. Participants had 1.5 s to indicate a guess
1324 (during this time a “?” was presented), while the subsequent feedback was presented for 1.0 s. In
1325 addition, there was a 1.0 s intertrial interval with a “+” on the screen. The task was presented in
1326 blocks that each included eight trials that were either mostly reward (6 reward trials that were
1327 pseudo-randomly interleaved with either 1 neutral and 1 loss trial, 2 neutral trials, or 2 loss trials)
1328 or mostly loss (6 loss trials interleaved with either 1 neutral and 1 reward trial, 2 neutral trials,
1329 or 2 reward trials) trials. In each of the two fMRI runs there were 2 mostly reward and 2 mostly
1330 loss blocks, interleaved with 4 fixation blocks (15 s each, during which a “+” is presented on the
1331 screen). All participants were provided with money as a result of completing the experiment. The
1332 amount they received was standardized due to the fixed nature of the experiment.

1333 **Motor** Participants were presented with visual cues that asked them to tap their left or right
1334 fingers, squeeze their left or right toes, or move their tongue. The task was presented in blocks of
1335 12 s that each included only one movement type (10 movements). Each block was preceded by a 3
1336 s cue. In each of the two fMRI runs, 13 blocks were presented with 2 blocks for tongue movements,
1337 4 blocks for hand movements (2 left, 2 right), and 4 blocks for foot movements (again, 2 left and
1338 2 right). In addition, three 15 s fixation blocks were included in each run.

1339 **Language** This task consisted of two runs that each interleaved 4 blocks of a story task and 4
1340 blocks of a math task. In the story task, participants were presented with brief auditory stories (5-9
1341 sentences) that were adapted from Aesop’s fables. After each story, a 2-alternative forced-choice
1342 question asked the participant about the topic of the story. In the math task, participants were
1343 similarly presented with an auditory math problem that asked them to complete 2-alternative
1344 forced choice addition or subtraction problems. For example, participants heard the operation
1345 “fourteen plus twelve”, followed by “equals” and then two choice alternatives (“twenty-nine or
1346 twenty-six”). Participants indicated with a button press whether they choose the first or second
1347 answer. The lengths of the blocks varied (with an average of approximately 30 s per block), but

1348 the task was designed in such a way that the math task blocks matched the length of the story task
1349 blocks (with some additional math trials at the end of a block if needed to complete the 3.8min
1350 run).

1351 **Social** Participants were presented with video clips (20 s each) that showed objects (squares, cir-
1352 cles, triangles) that either interacted in some way or were moving randomly. After each video clip,
1353 participants indicated whether they think that the objects had a social interaction (an interaction
1354 that appears as if the objects are taking into account each other’s feelings or thoughts), they are
1355 not sure, or they think the objects did not interact. Each of the two fMRI runs included 5 video
1356 blocks (2 with interaction and 3 without in one run and 3 with interaction and 2 without in the
1357 other run) as well as 5 15 s fixation blocks.

1358 **Relational** In this task, participants saw stimuli that were composed of six different shapes that
1359 were filled with one of six different textures. In the relational task condition, 2 pairs of objects
1360 were presented, one at the top of the screen and the other at the bottom. Participants were
1361 told that they should first decide what dimension (shape or texture) differs across the top pair
1362 of objects and then whether the bottom pair of objects differs along the same dimension. In the
1363 matching condition, participants were shown two objects at the top of the screen and one at the
1364 bottom. A word in the middle of the screen then indicated whether participants should decide if
1365 the bottom object matched either of the two top objects on the "shape" or "texture" dimension.
1366 In the relational condition, stimuli were presented for 3500 ms, with a 500 ms intertrial interval
1367 and four trials per block. In the matching condition, stimuli were presented for 2800 ms, with a
1368 400 ms intertrial interval, and a total of five trials per block. Each block lasted a total of 18 s. In
1369 each of the two fMRI runs three relational blocks, three matching blocks and three fixation blocks
1370 (16 s each) were presented.

1371 **Emotion** In emotion trials, participants were presented with with two faces at the bottom of
1372 the screen and one face at the top. These faces had an either angry or fearful expression. The
1373 participants were asked to decide which of the two faces on the bottom matches the face at the
1374 top. In neutral trials, participants were asked to decide which of two shapes at the bottom of the
1375 screen matches a shape that is presented at the top. In this task, trials were presented in blocks
1376 of six trials of the same task (face or shape). In each trial, the stimulus was presented for 2 s in
1377 addition to a 1 s intertrial interval. Each block was further preceded by a 3 s cue for the task
1378 (shape or face). Each of the two fMRI runs included three face and three shape blocks. Due to a
1379 bug in the experiment script, the experiment stopped before the final three trials of the last block
1380 of each trial (for further details on this bug, see Barch et al. (2013)).

1381 A.2 GLM analysis details

1382 **FMRI** Our GLM subject-level analyses of the fMRI data included one predictor for each of the
1383 four cognitive states in the design matrix (each representing a box-car function for the occurrence of
1384 a cognitive state). We convolved these predictors with a canonical glover haemodynamic response
1385 function (HRF; Lindquist et al., 2009) as implemented in Nilearn 0.8.0 (Abraham et al., 2014), to
1386 generate the model predictors. We added temporal derivative terms derived from each predictor,
1387 an intercept and an indicator of the experiment run to the design matrix, which we all treated as
1388 confounds of no interest. The derivative terms were computed by the use of the cosine drift model
1389 as implemented in Nilearn 0.8.0 (Abraham et al., 2014).

1390 To generate a set of group-level brain maps with the GLM, we computed a second-level GLM
1391 contrast by the use of the standard two-stage procedure for a random-effects group-level analysis,
1392 as proposed by Holmes and Friston (1998). Here, the subject-level regression coefficients β are

1393 treated as random effects in a second-level linear contrast analysis, where the distribution of first-
1394 level β -contrasts is assessed. Contrasts were computed between each cognitive state and all others.
1395 The resulting group-level brain maps show the z-scores resulting from this test.

1396 **Relevances** Our GLM analyses of the relevance data resulting from the application of the LRP
1397 technique to DeepLight’s decoding decisions (for an overview of the LRP technique, see section
1398 4.4 of the main text) included one predictor for each of the four cognitive states in the data (each
1399 representing a box-car function for the occurrence of a cognitive state). Our previous analyses
1400 have indicated that DeepLight’s relevance data show a similar temporal evolution as the HRF (see
1401 Fig. 6 of Thomas et al., 2019a). For this reason, we next convolved the predictors with a canonical
1402 glover HRF (Lindquist et al., 2009), as implemented in Nilearn 0.8.0 (Abraham et al., 2014), to
1403 generate a set of model predictors.

1404 We further added temporal derivative terms derived from each predictor, an intercept and
1405 an indicator of the experiment run to the design matrix. The temporal derivative terms were
1406 computed by the use of the cosine drift model as implemented in Nilearn 0.8.0 (Abraham et al.,
1407 2014). Additionally, we added one regressor to the design matrix indicating the total sum of
1408 relevance values contained in each fMRI volume (i.e., TR), to account for the variability in the
1409 sum of relevance values between TRs resulting from variability in the certainty of DeepLight’s
1410 predictions (for an overview of the LRP technique, see section 4.4 of the main text). To also
1411 account for non-linear relationships between this regressor and the relevance values, we added
1412 regressors for the first derivative of the relevance sums, the squared relevance sums, and the first
1413 derivative of the squared relevance sums to the design matrix. All of these predictors were treated
1414 as confounds of no interest.

1415 Lastly, we added two regressors to the design matrix indicating whether DeepLight correctly
1416 or incorrectly identified the cognitive state of each TR (again in form of two box-car functions).
1417 Importantly, we included these two predictors in each computed contrast, by contrasting each
1418 cognitive state against all other states and by contrasting correct versus incorrect predictions (e.g.,
1419 to compute a contrast for the body state in the HCP-WM task (see section 4.1.1 of the main text),
1420 we would set the contrast vector to: 3, -1, -1, -1, 1, -1 for the predictors: body, face, place, tool,
1421 correct, incorrect).

1422 To generate a set of group-level brain maps with the GLM, we computed a second-level GLM
1423 contrast by the use of the standard two-stage procedure for a random-effects group-level analysis,
1424 as proposed by Holmes and Friston (1998). Here, the subject-level regression coefficients β are
1425 treated as random effects in a second-level linear contrast analysis, where the distribution of first-
1426 level β -contrasts is assessed. The resulting group-level brain maps show the Z-values resulting
1427 from this test.

1428 **A.3 fMRIPrep details for Multi-task data**

1429 This dataset was processed using *fMRIPrep* 20.0.5 (Esteban et al. (2019); Esteban et al. (2018);
1430 RRID:SCR.016216), which is based on *Nipype* 1.4.2 (Gorgolewski et al. (2011); Gorgolewski et al.
1431 (2018); RRID:SCR.002502).

1432 **Anatomical data preprocessing** The T1-weighted (T1w) image was corrected for intensity
1433 non-uniformity (INU) with `N4BiasFieldCorrection` (Tustison et al., 2010), distributed with
1434 ANTs 2.2.0 (Avants et al., 2008, RRID:SCR.004757), and used as T1w-reference through-
1435 out the workflow. The T1w-reference was then skull-stripped with a *Nipype* implementa-
1436 tion of the `antsBrainExtraction.sh` workflow (from ANTs), using OASIS30ANTs as tar-
1437 get template. Brain tissue segmentation of cerebrospinal fluid (CSF), white-matter (WM)
1438 and gray-matter (GM) was performed on the brain-extracted T1w using `fast` (FSL 5.0.9,
1439 RRID:SCR.002823, Zhang et al., 2001). Volume-based spatial normalization to two standard

spaces (MNI152NLin6Asym, MNI152NLin2009cAsym) was performed through nonlinear registration with `antsRegistration` (ANTs 2.2.0), using brain-extracted versions of both T1w reference and the T1w template. The following templates were selected for spatial normalization: *FSL's MNI ICBM 152 non-linear 6th Generation Asymmetric Average Brain Stereotaxic Registration Model* [Evans et al. (2012), RRID:SCR_002823; TemplateFlow ID: MNI152NLin6Asym], *ICBM 152 Nonlinear Asymmetrical template version 2009c* [Fonov et al. (2009), RRID:SCR_008796; TemplateFlow ID: MNI152NLin2009cAsym],

Functional data preprocessing For each of the 18 BOLD runs found per subject (across all tasks and sessions), the following preprocessing was performed. First, a reference volume and its skull-stripped version were generated using a custom methodology of *fMRIPrep*. Susceptibility distortion correction (SDC) was omitted. The BOLD reference was then co-registered to the T1w reference using `flirt` (FSL 5.0.9, Jenkinson and Smith, 2001) with the boundary-based registration (Greve and Fischl, 2009) cost-function. Co-registration was configured with nine degrees of freedom to account for distortions remaining in the BOLD reference. Head-motion parameters with respect to the BOLD reference (transformation matrices, and six corresponding rotation and translation parameters) are estimated before any spatiotemporal filtering using `mcflirt` (FSL 5.0.9, Jenkinson et al., 2002). BOLD runs were slice-time corrected using `3dTshift` from AFNI 20160207 (Cox and Hyde, 1997, RRID:SCR_005927). The BOLD time-series (including slice-timing correction when applied) were resampled onto their original, native space by applying the transforms to correct for head-motion. These resampled BOLD time-series will be referred to as *preprocessed BOLD in original space*, or just *preprocessed BOLD*. The BOLD time-series were resampled into standard space, generating a *preprocessed BOLD run in MNI152NLin6Asym space*. First, a reference volume and its skull-stripped version were generated using a custom methodology of *fMRIPrep*. Several confounding time-series were calculated based on the *preprocessed BOLD*: framewise displacement (FD), DVARS and three region-wise global signals. FD and DVARS are calculated for each functional run, both using their implementations in *Nipype* (following the definitions by Power et al., 2014). The three global signals are extracted within the CSF, the WM, and the whole-brain masks. Additionally, a set of physiological regressors were extracted to allow for component-based noise correction (*CompCor*, Behzadi et al., 2007). Principal components are estimated after high-pass filtering the *preprocessed BOLD* time-series (using a discrete cosine filter with 128s cut-off) for the two *CompCor* variants: temporal (tCompCor) and anatomical (aCompCor). tCompCor components are then calculated from the top 5% variable voxels within a mask covering the subcortical regions. This subcortical mask is obtained by heavily eroding the brain mask, which ensures it does not include cortical GM regions. For aCompCor, components are calculated within the intersection of the aforementioned mask and the union of CSF and WM masks calculated in T1w space, after their projection to the native space of each functional run (using the inverse BOLD-to-T1w transformation). Components are also calculated separately within the WM and CSF masks. For each CompCor decomposition, the k components with the largest singular values are retained, such that the retained components' time series are sufficient to explain 50 percent of variance across the nuisance mask (CSF, WM, combined, or temporal). The remaining components are dropped from consideration. The head-motion estimates calculated in the correction step were also placed within the corresponding confounds file. The confound time series derived from head motion estimates and global signals were expanded with the inclusion of temporal derivatives and quadratic terms for each (Satterthwaite et al., 2013). Frames that exceeded a threshold of 0.5 mm FD or 1.5 standardised DVARS were annotated as motion outliers. All resamplings can be performed with a *single interpolation step* by composing all the pertinent transformations (i.e. head-motion transform matrices, susceptibility distortion correction when available, and co-registrations to anatomical and output spaces). Gridded

1490 (volumetric) resamplings were performed using `antsApplyTransforms` (ANTs), configured
1491 with Lanczos interpolation to minimize the smoothing effects of other kernels (Lanczos, 1964).
1492 Non-gridded (surface) resamplings were performed using `mri_vol2surf` (FreeSurfer).

1493 Many internal operations of *fMRIPrep* use *Nilearn* 0.6.2 (Abraham et al., 2014), mostly within
1494 the functional processing workflow. For more details of the pipeline, see the section corresponding
1495 to workflows in *fMRIPrep*'s documentation.

1496 The above boilerplate text was automatically generated by *fMRIPrep* with the express intention
1497 that users should copy and paste this text into their manuscripts *unchanged*. It is released under
1498 the CC0 license.

1499 A.4 fMRIPrep details for HCP working memory task

1500 Results included in this manuscript come from preprocessing performed using *fMRIPrep* 20.0.5
1501 (Esteban et al. (2019); Esteban et al. (2018); RRID:SCR_016216), which is based on *Nipype* 1.4.2
1502 (Gorgolewski et al. (2011); Gorgolewski et al. (2018); RRID:SCR_002502).

1503 **Anatomical data preprocessing** The T1-weighted (T1w) image was corrected for intensity
1504 non-uniformity (INU) with `N4BiasFieldCorrection` (Tustison et al., 2010), distributed with
1505 ANTs 2.2.0 (Avants et al., 2008, RRID:SCR_004757), and used as T1w-reference through-
1506 out the workflow. The T1w-reference was then skull-stripped with a *Nipype* implementa-
1507 tion of the `antsBrainExtraction.sh` workflow (from ANTs), using OASIS30ANTs as tar-
1508 get template. Brain tissue segmentation of cerebrospinal fluid (CSF), white-matter (WM)
1509 and gray-matter (GM) was performed on the brain-extracted T1w using `fast` (FSL 5.0.9,
1510 RRID:SCR_002823, Zhang et al., 2001). Volume-based spatial normalization to two standard
1511 spaces (MNI152NLin6Asym, MNI152NLin2009cAsym) was performed through nonlinear reg-
1512 istration with `antsRegistration` (ANTs 2.2.0), using brain-extracted versions of both T1w
1513 reference and the T1w template. The following templates were selected for spatial nor-
1514 malization: *FSL's MNI ICBM 152 non-linear 6th Generation Asymmetric Average Brain*
1515 *Stereotaxic Registration Model* [Evans et al. (2012), RRID:SCR_002823; TemplateFlow ID:
1516 MNI152NLin6Asym], *ICBM 152 Nonlinear Asymmetrical template version 2009c* [Fonov
1517 et al. (2009), RRID:SCR_008796; TemplateFlow ID: MNI152NLin2009cAsym],

1518 **Functional data preprocessing** For each of the 14 BOLD runs found per subject (across all
1519 tasks and sessions), the following preprocessing was performed. First, a reference volume and
1520 its skull-stripped version were generated using a custom methodology of *fMRIPrep*. Suscep-
1521 tibility distortion correction (SDC) was omitted. The BOLD reference was then co-registered
1522 to the T1w reference using `flirt` (FSL 5.0.9, Jenkinson and Smith, 2001) with the boundary-
1523 based registration (Greve and Fischl, 2009) cost-function. Co-registration was configured
1524 with nine degrees of freedom to account for distortions remaining in the BOLD reference.
1525 Head-motion parameters with respect to the BOLD reference (transformation matrices, and
1526 six corresponding rotation and translation parameters) are estimated before any spatiotem-
1527 poral filtering using `mcflirt` (FSL 5.0.9, Jenkinson et al., 2002). The BOLD time-series
1528 (including slice-timing correction when applied) were resampled onto their original, native
1529 space by applying the transforms to correct for head-motion. These resampled BOLD time-
1530 series will be referred to as *preprocessed BOLD in original space*, or just *preprocessed BOLD*.
1531 The BOLD time-series were resampled into standard space, generating a *preprocessed BOLD*
1532 *run in MNI152NLin6Asym space*. First, a reference volume and its skull-stripped version
1533 were generated using a custom methodology of *fMRIPrep*. Several confounding time-series
1534 were calculated based on the *preprocessed BOLD*: framewise displacement (FD), DVARS
1535 and three region-wise global signals. FD and DVARS are calculated for each functional run,
1536 both using their implementations in *Nipype* (following the definitions by Power et al., 2014).

1537 The three global signals are extracted within the CSF, the WM, and the whole-brain masks.
1538 Additionally, a set of physiological regressors were extracted to allow for component-based
1539 noise correction (*CompCor*, Behzadi et al., 2007). Principal components are estimated after
1540 high-pass filtering the *preprocessed BOLD* time-series (using a discrete cosine filter with 128s
1541 cut-off) for the two *CompCor* variants: temporal (tCompCor) and anatomical (aCompCor).
1542 tCompCor components are then calculated from the top 5% variable voxels within a mask
1543 covering the subcortical regions. This subcortical mask is obtained by heavily eroding the
1544 brain mask, which ensures it does not include cortical GM regions. For aCompCor, compo-
1545 nents are calculated within the intersection of the aforementioned mask and the union
1546 of CSF and WM masks calculated in T1w space, after their projection to the native space
1547 of each functional run (using the inverse BOLD-to-T1w transformation). Components are
1548 also calculated separately within the WM and CSF masks. For each *CompCor* decomposi-
1549 tion, the k components with the largest singular values are retained, such that the retained
1550 components' time series are sufficient to explain 50 percent of variance across the nuisance
1551 mask (CSF, WM, combined, or temporal). The remaining components are dropped from
1552 consideration. The head-motion estimates calculated in the correction step were also placed
1553 within the corresponding confounds file. The confound time series derived from head motion
1554 estimates and global signals were expanded with the inclusion of temporal derivatives and
1555 quadratic terms for each (Satterthwaite et al., 2013). Frames that exceeded a threshold of 0.5
1556 mm FD or 1.5 standardised DVARS were annotated as motion outliers. All resamplings can
1557 be performed with *a single interpolation step* by composing all the pertinent transformations
1558 (i.e. head-motion transform matrices, susceptibility distortion correction when available, and
1559 co-registrations to anatomical and output spaces). Gridded (volumetric) resamplings were
1560 performed using `antsApplyTransforms` (ANTs), configured with Lanczos interpolation to
1561 minimize the smoothing effects of other kernels (Lanczos, 1964). Non-gridded (surface) re-
1562 samplings were performed using `mri_vol2surf` (FreeSurfer).

1563 Many internal operations of *fMRIPrep* use *Nilearn* 0.6.2 (Abraham et al., 2014), mostly within
1564 the functional processing workflow. For more details of the pipeline, see the section corresponding
1565 to workflows in *fMRIPrep*'s documentation.

1566 **Copyright Waiver** The above boilerplate text was automatically generated by *fMRIPrep* with
1567 the express intention that users should copy and paste this text into their manuscripts *unchanged*.
1568 It is released under the CC0 license

B Supplement to results

B.1 Do basic statistical differences between HCP and Multi-task data affect transfer performance?

To better understand whether any basic differences in statistical properties, noise, or preprocessing between the HCP and Multi-task datasets affected the transfer performance of the pre-trained 3D-DeepLight variant, we performed a sequence of additional analyses.

We can immediately rule out basic differences in the temporal distribution of the voxel signals as we detrended and standardized the time series signal of each voxel within each fMRI run (to have a mean of 0 and unit variance; see section 4.1 of the main text). DeepLight further does not know about the temporal distribution of brain activity as it solely acts on the level of individual fMRI volumes. We therefore next probed the mean and standard deviation of voxel activities within each fMRI volume. We did not find any meaningful differences in the distribution of the volume means and standard deviations between the HCP and Multi-task datasets (see Appendix Fig. B.2).

We also tested whether other generic differences in noise between the HCP and Multi-task datasets affected transfer performance, by performing a confound correction of the Multi-task fMRI data, in which we regressed out variance related to the six motion correction parameters and three temporal and anatomical noise components resulting from fMRIPrep’s CompCor method (for an overview, see Appendix Fig. B.4). Yet, the pre-trained model did not perform better when fine-tuned on the confound-corrected fMRI data than when fine-tuned on the fMRI data that was not confound corrected (for an overview of the training methods, see section 4.3 of the main text). 3D-DeepLight’s final decoding accuracy on the confound-corrected data was 43.27%, thereby -2.5% worse than when applied to the uncorrected fMRI data ($t(5) = -4.65, P = 0.0056$; Appendix Fig. B.4).

Lastly, we also tested whether the transfer of the pre-trained model to the Multi-task data was affected by the different preprocessing that we applied to both datasets (we preprocessed the Multi-task dataset with fMRIPrep (Esteban et al., 2019), whereas the HCP uses an internal preprocessing pipeline; see section 4.1 of the main text). To this end, we downloaded the raw fMRI data of another 50 subjects in the HCP working memory task and also preprocessed these with fMRIPrep (for an overview of the preprocessing steps, see Appendix A.4). Interestingly, the pre-trained 3D-DeepLight variant again exhibited the advantages of transfer learning in this newly preprocessed fMRI dataset, by learning faster and achieving higher decoding accuracies than a model variant that was not pre-trained (see Appendix Fig. B.3). After training on the fMRI data of 20 subjects from this newly preprocessed dataset, the pre-trained model achieved a final decoding accuracy of 72.95% in the fMRI data of the remaining 30 subjects, while the model variant that was not pre-trained achieved a final decoding accuracy of 64.46% (i.e., -8.49% worse than the pre-trained model, $t(29) = -13.28, P < 0.0001$; see Appendix Fig. B.3; for an overview of the training methods, see section 4.3 of the main text).

Overall, we can therefore rule out that the transfer of the pre-trained model to the Multi-task dataset was affected by basic differences in the statistical properties, noise or preprocessing between the HCP and Multi-task datasets.

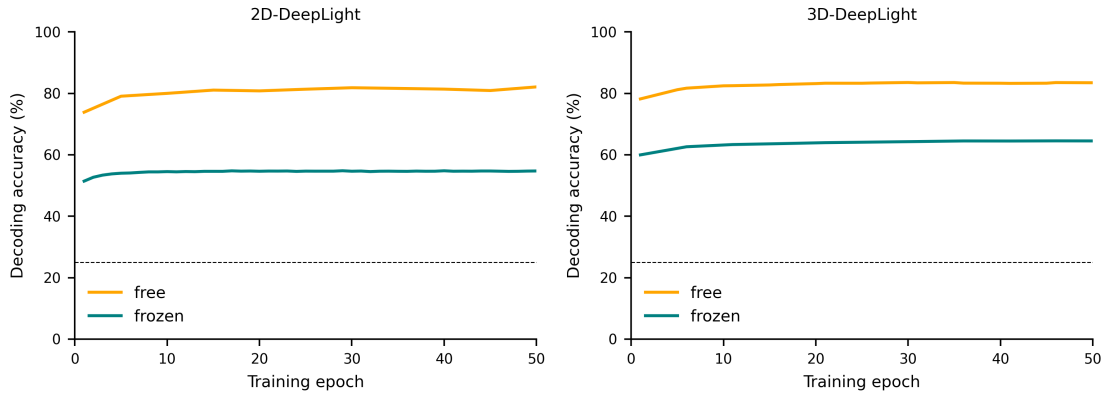


Figure B.1: Comparing two different fine-tuning approaches on the validation data of the HCP working memory task (see section 4.1.1 of the main text). We initialized the weights of two variants of each DeepLight architecture (left: 2D-DeepLight, right: 3D-DeepLight) to the weights of the pre-trained models (all except for the output layer, which now included four instead of 16 neurons; see section 4.2 of the main text for an overview of the architectures and Fig. 3 of the main text for an overview of the pre-trained model performance). We froze the pre-trained weights of one variant of each architecture during fine-tuning (depicted in green), while the other model variant was allowed to train all of its weights during fine-tuning (depicted in yellow) (see section 4.3 of the main text for an overview of the training procedures). Lines indicate decoding accuracy in the validation data as a function of the training epochs. Chance accuracy is indicated by the dashed horizontal line.

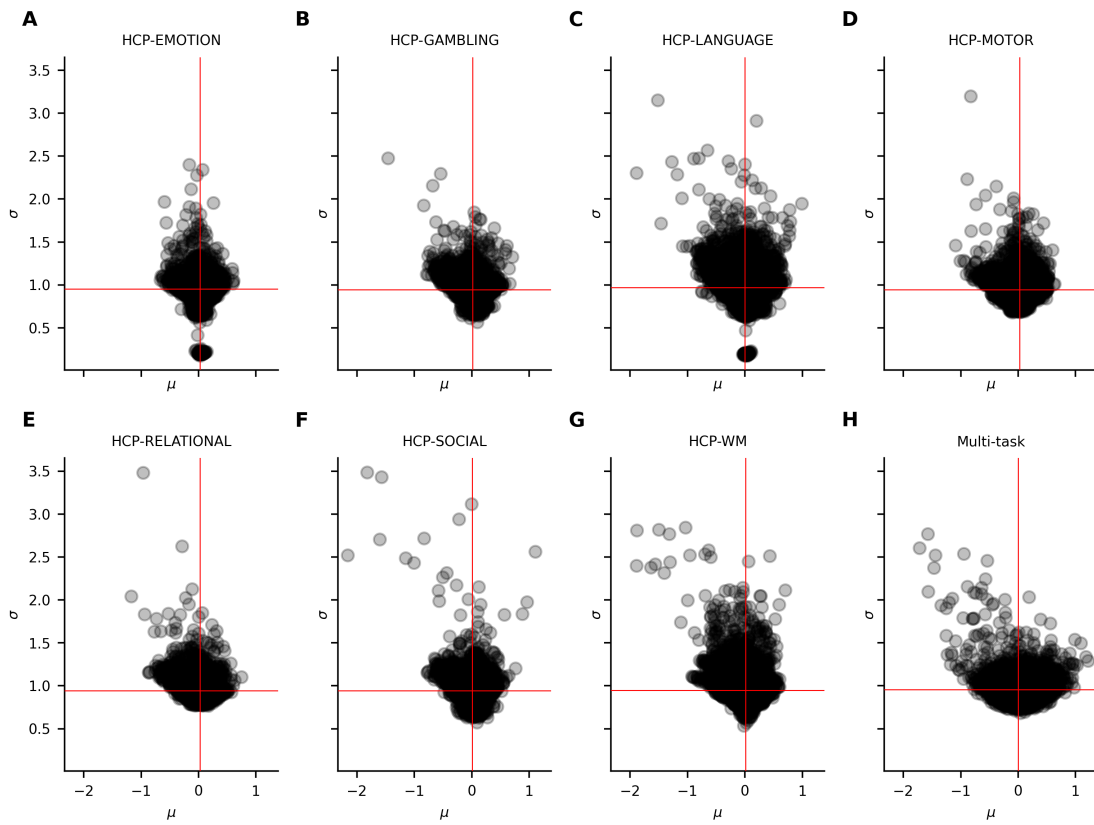


Figure B.2: Mean and standard deviation of voxel activities within each preprocessed fMRI volume in the validation datasets of the HCP experimental tasks (A-G) and Multi-task data (H) (for an overview of the datasets, see section 4.1 of the main text). Scatter points indicate individual fMRI volumes. Red lines indicate the mean over volumes.

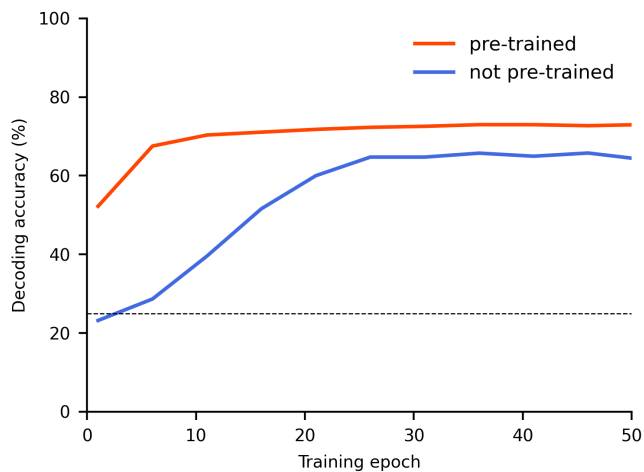


Figure B.3: Training decoding accuracy for a pre-trained (red) and not pre-trained (blue) 3D-DeepLight variant in the validation data of the HCP working memory task that was preprocessed with fMRIPrep (see Appendix B.1; see section 4.3 of the main text for an overview of the training procedures). An epoch was defined as an entire iteration over the training dataset. Lines indicate decoding accuracy. Chance accuracy is indicated by the dashed horizontal line.

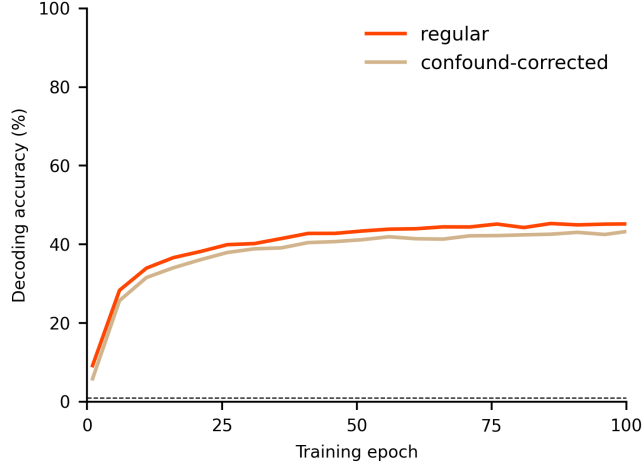


Figure B.4: Training decoding accuracy for the pre-trained 3D-DeepLight variant in two conditions: when it is fine-tuned on the regular fMRI data of the Multi-task dataset (red) or on a version that is corrected for basic noise confounds (tan). Specifically, we corrected the Mutli-task data for any variance resulting from the six parameters of basic motion correction, as well as the three temporal and anatomical noise components with the largest singular values resulting from fMRIPrep’s CompCor method (for details on this method, see Behzadi et al. (2007)), by regressing their variance out of the time-series signal of each voxel (as implemented in Nilearns ”signal.clean” function; Abraham et al., 2014). See section 4.3 of the main text for an overview of the training procedures. An epoch was defined as an entire iteration over the training dataset. Lines indicate decoding accuracy. Chance accuracy is indicated by the dashed horizontal line.

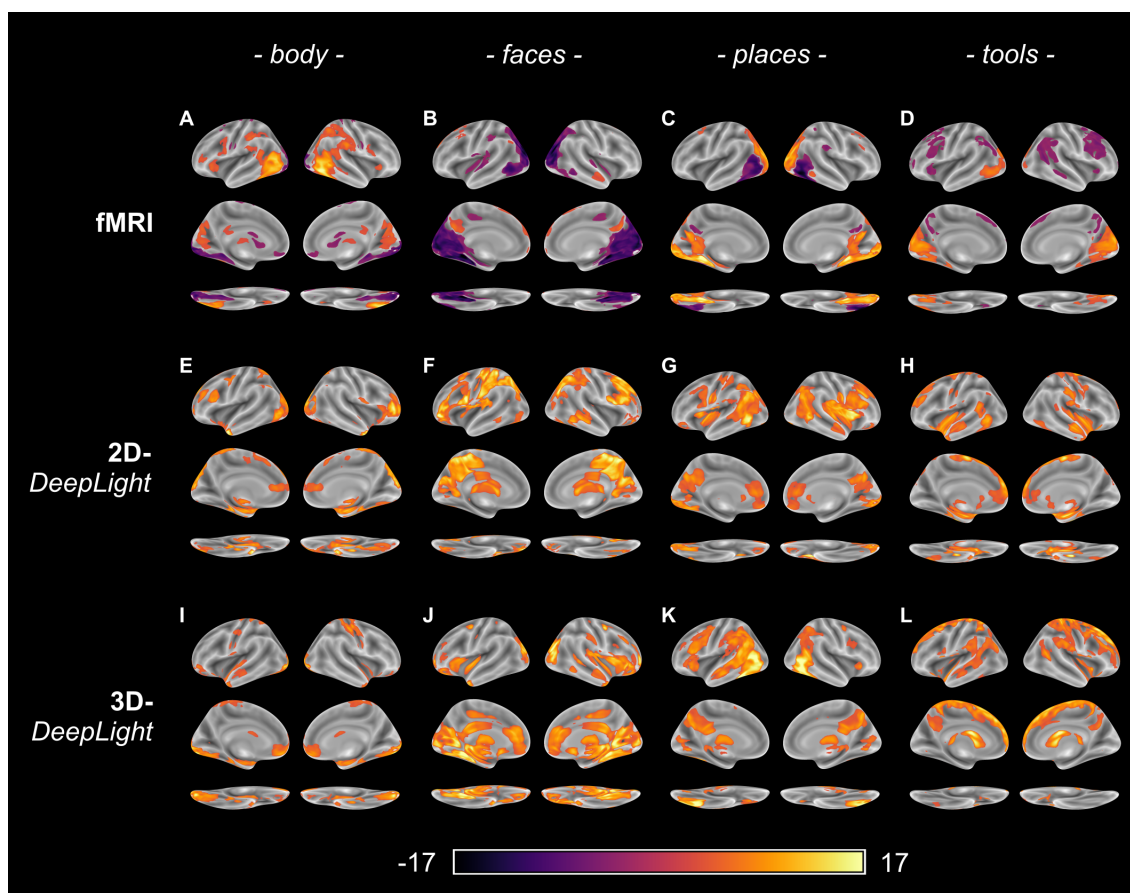


Figure B.5: Learned mappings between brain activity and cognitive states of the pre-trained DeepLight variants that were fine-tuned on the full training dataset of the HCP-WM experimental task (see section 2.3 of the main text). A-D: We first computed a standard two-stage GLM analysis (Holmes and Friston, 1998) of the fMRI data of the 50 subjects in the validation dataset of this task. E-L: We then also interpreted the decoding decisions of the 2D- (E-H) and 3D-DeepLight (I-L) variants for the same data. To identify the brain regions that each DeepLight variant associates most strongly with a cognitive state, we computed a similar two-stage GLM analysis of the resulting relevance data (restricting the resulting z-scores to only positive values). All GLM analyses were performed on parcellated brain data by the use of the dictionaries for functional modes (DiFuMo) atlas with 256 brain networks (Dadi et al., 2020) and computed separately for each experimental task by contrasting each cognitive state of the task against all other states of that task (for details on the GLM analysis, see Appendix A.2). All brain maps are thresholded at a false-discovery rate of 0.001 and projected onto the inflated cortical surface of the FsAverage template (Fischl, 2012). Brighter yellow values indicate larger z-scores.

Declaration of interests

¹⁶¹⁰

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: