

# 1 Analyzing Neuroimaging Data Through Recurrent Deep 2 Learning Models

3 **Armin W. Thomas**<sup>1,2,3,4</sup>, **Hauke R. Heekeren**<sup>2,4 \*</sup>, **Klaus-Robert Müller**<sup>1,5,6 \*</sup>,  
4 **Wojciech Samek**<sup>7\*</sup>

5 1 Machine Learning Group, Technische Universität Berlin, Berlin, Germany

6 2 Center for Cognitive Neuroscience Berlin, Freie Universität Berlin, Berlin, Germany

7 3 Max Planck School of Cognition, Leipzig, Germany

8 4 Department of Education & Psychology, Freie Universität Berlin, Berlin, Germany

9 5 Department of Brain & Cognitive Engineering, Korea University, Seoul, South Korea

10 6 Max Planck Institute for Informatics, Saarbrücken, Germany

11 7 Machine Learning Group, Fraunhofer Heinrich Hertz Institute, Berlin, Germany

12 \* **Correspondence:**

13 hauke.heekeren@fu-berlin.de, klaus-robot.mueller@tu-berlin.de,  
14 wojciech.samek@hhi.fraunhofer.de

**Keywords:** decoding, neuroimaging, fMRI, whole-brain, deep learning, recurrent, interpretability

17 Manuscript length: 8900 words

18    Figures: 6 (main) + 5 (supplement)

## 19 Abstract

20 The application of deep learning (DL) models to neuroimaging data poses several  
21 challenges, due to the high dimensionality, low sample size and complex temporo-  
22 spatial dependency structure of these data. Even further, DL models often act as as  
23 *black boxes*, impeding insight into the association of cognitive state and brain activity.  
24 To approach these challenges, we introduce the DeepLight framework, which utilizes  
25 long short-term memory (LSTM) based DL models to analyze *whole-brain* functional  
26 Magnetic Resonance Imaging (fMRI) data. To decode a cognitive state (e.g., seeing the  
27 image of a house), DeepLight separates an fMRI volume into a sequence of axial brain  
28 slices, which is then sequentially processed by an LSTM. To maintain interpretability,  
29 DeepLight adapts the layer-wise relevance propagation (LRP) technique. Thereby,  
30 decomposing its decoding decision into the contributions of the single input voxels to  
31 this decision. Importantly, the decomposition is performed on the level of single fMRI  
32 volumes, enabling DeepLight to study the associations between cognitive state and  
33 brain activity on several levels of data granularity, from the level of the group down to  
34 the level of single time points. To demonstrate the versatility of DeepLight, we apply it  
35 to a large fMRI dataset of the Human Connectome Project. We show that DeepLight  
36 outperforms conventional approaches of uni- and multivariate fMRI analysis in  
37 decoding the cognitive states and in identifying the physiologically appropriate brain  
38 regions associated with these states. We further demonstrate DeepLight’s ability to  
39 study the fine-grained temporo-spatial variability of brain activity over sequences of  
40 single fMRI samples.

## 1. Introduction

Neuroimaging research has recently started collecting large corpora of experimental functional Magnetic Resonance Imaging (fMRI) data, often comprising many hundred individuals (e.g., Poldrack et al., 2013; Van Essen et al., 2013). By collecting these datasets, researchers want to gain insights into the associations between the cognitive states of an individual (e.g., while viewing images or performing a specific task) and the underlying brain activity, while also studying the variability of these associations across the population.

At first sight, the analysis of neuroimaging data thereby seems ideally suited for the application of deep learning (DL; Goodfellow et al., 2016; LeCun et al., 2015) methods, due to the availability of large and structured datasets. Generally, DL can be described as a class of representation-learning methods, with multiple levels of abstraction. At each level, the representation of the input data is transformed by a simple, but non-linear function. The resulting hierarchical structure of non-linear transforms enables DL methods to learn complex functions. It also enables them to identify intricate signals in noisy data, by projecting the input data into a higher-level representation, in which those aspects of the input data that are irrelevant to identify an analysis target are suppressed and those that are relevant are amplified. With this higher-level perspective, DL methods can associate a target variable with variable patterns in the input data. Importantly, DL methods can autonomously learn these projections from the data and therefore do not require a thorough prior understanding of the mapping between input data and analysis target (for a detailed discussion, see LeCun et al., 2015). For these reasons, DL methods seem ideally suited for the analysis of neuroimaging data, where intricate, highly variable patterns of brain activity are hidden in large, high-dimensional datasets and the mapping between cognitive state and brain activity is often unknown.

While researchers have started exploring the application of DL models to neuroimaging data (e.g., Mensch et al., 2018; Nie et al., 2016; Petrov et al., 2018; Plis et al., 2014; Sarraf and Tofighi, 2016; Suk et al., 2014; Yousefnezhad and Zhang, 2018), two major challenges have so far prevented broad DL usage: (1) Neuroimaging data are high dimensional, while containing comparably few samples. For example, a typical fMRI dataset comprises up to a few hundred samples per subject and recently up to several hundred subjects (e.g., Van Essen et al., 2013), while each sample contains several hundred thousand dimensions (i.e., voxels). In such analysis settings, DL models (as well as more traditional machine learning approaches) are likely to suffer from overfitting (by too closely capturing those dynamics that are specific to the training data, so that their predictive performance does not generalize well to new data). (2) DL models have often been considered as non-linear *black box models*, disguising the relationship between input data and decoding decision. Thereby, impeding insight into (and interpretation of) the association between cognitive state and brain activity.

To approach these challenges, we propose the DeepLight framework, which defines a method to utilize long short-term memory (LSTM) based DL architectures (Donahue et al., 2015; Hochreiter and Schmidhuber, 1997) to analyze whole-brain neuroimaging data. In DeepLight, each whole-brain volume is sliced into a sequence of axial images. To decode an underlying cognitive state, the resulting sequence of images is processed

by a combination of convolutional and recurrent DL elements. Thereby, DeepLight successfully copes with the high dimensionality of neuroimaging data, while modeling the full spatial dependency structure of whole-brain activity (within and across axial brain slices). Conceptually, DeepLight builds upon the searchlight approach. Instead of moving a small searchlight beam around in space, DeepLight explores brain activity more in-depth, by looking through the full sequence of axial brain slices, before making a decoding decision. To subsequently relate brain activity and cognitive state, DeepLight applies the layer-wise relevance propagation (LRP; Bach et al., 2015; Lapuschkin et al., 2016) method to its decoding decisions. Thereby, decomposing these decisions into the contributions of the single input voxels to each decision. Importantly, the LRP analysis is performed on the level of a single input samples, enabling an analysis on several levels of data granularity, from the level of the group down to the level of single subjects, trials and time points. These characteristics make DeepLight ideally suited to study the fine-grained temporo-spatial distribution of brain activity underlying sequences of single fMRI samples.

Here, we will demonstrate the versatility of DeepLight, by applying it to an openly available fMRI dataset of the Human Connectome Project (Van Essen et al., 2013). In particular, to the data of an N-back task, in which 100 subjects viewed images of either body parts, faces, places or tools in two separate fMRI experiment runs (for an overview, see Section 2.1 and Supplementary Fig. S1). Subsequently, we will evaluate the performance of DeepLight in decoding the four underlying cognitive states (resulting from viewing an image of either of the four stimulus classes) from the fMRI data and identifying the brain regions associated with these states. To this end, we will compare the performance of DeepLight to three representative conventional approaches to the uni- and multivariate analysis of neuroimaging data, with widespread application in the literature. In particular, we will compare DeepLight to the General Linear Model (GLM; Friston et al., 1994), searchlight analysis (Kriegeskorte et al., 2006) and whole-brain Least Absolute Shrinkage Logistic Regression (whole-brain Lasso; Großenick et al., 2013; Wager et al., 2013). Note that the four analysis approaches differ in the number of voxels they include in their analyses. While the GLM analyses the data of single voxels independent of one another (univariate), the searchlight analysis utilizes the data of clusters of multiple voxels (multivariate) and the whole-brain lasso utilizes the data of all voxels in the brain (whole-brain). In this comparison, we find that DeepLight (1) decodes the cognitive states underlying the fMRI data more accurately than these other approaches, (2) improves its decoding performance better with growing datasets, (3) accurately identifies the physiologically appropriate associations between cognitive states and brain activity and (4) identifies these associations on multiple levels of data granularity (namely, the level of the group, subject, trial and time point). We also demonstrate DeepLight’s ability to study the temporo-spatial distribution of brain activity over a sequence of single fMRI samples.

## 2. Methods

## 126 2.1 Experiment paradigm

127 100 participants performed a version of the N-back task in two separate fMRI runs (for  
 128 an overview, see Supplementary Fig. S1 and Barch et al., 2013). Each of the two runs  
 129 (260s each) consisted of eight task blocks (25s each) and four fixation blocks (15s  
 130 each). Within each run, the four different stimulus types (body, face, place and tool)  
 131 were presented in separate blocks. Half of the task blocks used a 2-back working  
 132 memory task (participants were asked to respond "target" when the current stimulus was  
 133 the same as the stimulus 2 back) and the other half a 0-back working memory task (a  
 134 target cue was presented at the beginning of each block and the participants were asked  
 135 to respond "target" whenever the target cue was presented in the block). Each task block  
 136 consisted of 10 trials (2.5s each). In each trial, a stimulus was presented for 2s followed  
 137 by a 500 ms interstimulus interval (ISI). We were not interested in identifying any effect  
 138 of the N-back task condition on the evoked brain activity and therefore pooled the data  
 139 of both N-back conditions.

## 140 2.2 FMRI data acquisition & preprocessing

141 Functional MRI data of 100 unrelated participants for this experiment were provided in  
 142 a preprocessed format by the Human Connectome Project (HCP S1200 release), WU  
 143 Minn Consortium (Principal Investigators: David Van Essen and Kamil Ugurbil;  
 144 1U54MH091657) funded by the 16 NIH Institutes and Centers that support the NIH  
 145 Blueprint for Neuroscience Research; and by the McDonnell Center for Systems  
 146 Neuroscience at Washington University. Whole-brain EPI acquisitions were acquired  
 147 with a 32 channel head coil on a modified 3T Siemens Skyra with TR=720 ms, TE=33.1  
 148 ms, flip angle=52 deg, BW=2290 Hz/Px, in-plane FOV=208 × 180 mm, 72 slices, 2.0  
 149 mm isotropic voxels with a multi-band acceleration factor of 8. Two runs were acquired,  
 150 one with a right-to-left and the other with a left-to-right phase encoding (for further  
 151 methodological details on fMRI data acquisition, see Ugurbil et al., 2013).

152 The Human Connectome Project preprocessing pipeline for functional MRI data  
 153 ("fMRIVolume"; Glasser et al., 2013) includes the following steps: gradient unwarping,  
 154 motion correction, fieldmap-based EPI distortion correction, brain-boundary based  
 155 registration of EPI to structural T1-weighted scan, non-linear registration into MNI152  
 156 space, and grand-mean intensity normalization (for further details, see Glasser et al.,  
 157 2013; Ugurbil et al., 2013). In addition to the minimal preprocessing of the fMRI data  
 158 that was performed by the Human Connectome Project, we applied the following  
 159 preprocessing steps to the data for all decoding analyses: volume-based smoothing of  
 160 the fMRI sequences with a 3mm Gaussian kernel, linear detrending and standardization  
 161 of the single voxel signal time-series (resulting in a zero-centered voxel time-series with  
 162 unit variance) and temporal filtering of the single voxel time-series with a butterworth  
 163 highpass filter and a cutoff of 128s, as implemented in Nilearn 0.4.1 (Abraham et al.,  
 164 2014). In line with previous work (Jang et al., 2017), we further applied an outer brain  
 165 mask to each fMRI volume. We first identified those voxels whose activity was larger  
 166 than 5% of the maximum voxel signal within the fMRI volume and then only kept those  
 167 voxels for further analysis that were positioned between the first and last voxel to fulfill  
 168 this property in the three spatial dimensions of any functional brain volume of our  
 169 dataset. This resulted in a brain mask spanning  $74 \times 92 \times 81$  voxels ( $X \times Y \times Z$ ).

170 All of our preprocessing was performed by the use of Nilearn 0.4.1 (Abraham et al.,  
 171 2014). Importantly, we did not exclude any TR of an experiment block of the four  
 172 stimulus classes from the decoding analyses. However, we removed all fixation blocks  
 173 from the decoding analyses. Lastly, we split the fMRI data of the 100 subjects contained  
 174 in the dataset into two distinct training and test datasets (each containing the data of 70  
 175 and 30 randomly assigned subjects). All analyses presented throughout the following  
 176 solely include the data of the 30 subjects contained in the held-out test dataset (if not  
 177 stated otherwise).

## 178 **2.3 Data availability**

179 The data that support the findings of this study are openly available at the  
 180 ConnectomeDB S1200 Project page of the Human Connectome Project  
 181 (<https://db.humanconnectome.org/data/projects/HCP1200>).

## 182 **2.4. Baseline methods**

### 183 **2.4.1 General linear model**

184 The General Linear Model (GLM; Friston et al., 1994) represents a univariate brain  
 185 encoding model (Kriegeskorte and Douglas, 2018; Naselaris et al., 2011). Its goal is to  
 186 identify an association between cognitive state and brain activity, by predicting the time  
 187 series signal of a voxel from a set of experiment predictor:

$$188 \quad Y = X\beta + \epsilon \quad (1)$$

189 Here,  $Y$  presents a  $T \times N$  dimensional matrix containing the multivariate time series data  
 190 of  $N$  voxels and  $T$  time points.  $X$  represents the design matrix, which is composed of  
 191  $T \times P$  data points, where each column represents one of  $P$  predictors. Typically, each  
 192 predictor represents a variable that is manipulated during the experiment (e.g., the  
 193 presentation times of stimuli of one of the four stimulus classes).  $\beta$  represents a  $P \times N$   
 194 dimensional matrix of regression coefficients. To mimic the blood-oxygen-level  
 195 dependent (BOLD) response measured by the fMRI, each predictor is first convolved  
 196 with a hemodynamic response function (HRF; Lindquist et al., 2009), before fitting the  
 197  $\beta$ -coefficients to the data. After fitting, the resulting brain map of  $\beta$ -coefficients  
 198 indicates the estimated contribution of each predictor to the time series signal of each of  
 199 the  $N$  voxels.  $\epsilon$  represents a  $T \times N$  dimensional matrix of error terms. Importantly, the  
 200 GLM analyzes the time series signal of each voxel independently and thereby includes a  
 201 separate set of regression coefficients for each voxel in the brain.

### 202 **2.4.2 Searchlight analysis**

203 The searchlight analysis (Kriegeskorte et al., 2006) is a multivariate brain decoding  
 204 model (Kriegeskorte and Douglas, 2018; Naselaris et al., 2011). Its goal is to identify an  
 205 association between cognitive state and brain activity, by probing the ability of a  
 206 statistical classifier to identify the cognitive state from the activity pattern of a small  
 207 clusters of voxels. To this end, the entire brain is scanned with a sphere of a given radius  
 208 (the searchlight) and the performance of the classifier in decoding the cognitive states is  
 209 evaluated at each location, resulting in a brain map of decoding accuracies. These

210 decoding accuracies indicate how much information about the cognitive state is  
 211 contained in the activity pattern of the underlying cluster of voxels. Here, we used a  
 212 searchlight radius of 5.6mm and a linear-kernel Support Vector Machine (SVM)  
 213 classifier (if not reported otherwise).

214 Given a training dataset of  $T$  data points  $[y_t, x_t]_{t=1}^T$ , where  $x_t$  represents the activity  
 215 pattern of a cluster of voxels at time point  $t$  and  $y_t$  the corresponding label, the SVM  
 216 (Cortes and Vapnik, 1995) is defined as follows:

$$217 \quad \hat{y}(x) = \text{sign} \left[ \sum_{t=1}^T \alpha_t y_t \gamma(x, x_t) + b \right] \quad (2)$$

218 Here,  $\alpha_t$  and  $b$  are positive constants, whereas  $\gamma(x, x_t)$  represents the kernel function.  
 219 We used a linear kernel function, as implemented in Nilearn 0.4.1 (Abraham et al.,  
 220 2014). We then defined the decoding accuracy achieved by the searchlight analysis as  
 221 the maximum decoding accuracy that was achieved at any searchlight location in the  
 222 brain. Similarly, we used the searchlight location that achieved the highest decoding  
 223 accuracy to make decoding predictions (for example, to compute the confusion matrix  
 224 presented in Fig. 2C).

### 225 2.4.3 Whole-brain Least Absolute Shrinkage Logistic Regression

226 The whole-brain Least Absolute Shrinkage Logistic Regression (or whole-brain lasso;  
 227 Grosenick et al., 2013; Wager et al., 2013) represents a whole-brain decoding model  
 228 (Kriegeskorte and Douglas, 2018; Naselaris et al., 2011). It identifies an association  
 229 between cognitive state and brain activity, by probing the ability of a logistic model to  
 230 decode the cognitive state from whole-brain activity (with one logistic coefficient  $\beta_i$  per  
 231 voxel  $i$  in the brain). To reduce the risk of overfitting, resulting from the large number  
 232 of model coefficients, the whole-brain lasso applies Least Absolute Shrinkage  
 233 regularization to the likelihood function of the logistic model (Tibshirani, 1996;  
 234 Tikhonov, 1943). Thereby, forcing the logistic model to perform automatic variable  
 235 selection during parameter estimation, resulting in sparse coefficient estimates (i.e., by  
 236 forcing many coefficient estimates to be exactly 0). In particular, the optimization  
 237 problem of the whole-brain lasso can be defined as follows (again,  $N$  denotes the  
 238 number of voxels in the brain,  $T$  the number of fMRI sampling time points and  $[y_t, x_t]_{t=1}^T$   
 239 the set of class labels and voxel values of each fMRI sample):

$$240 \quad \min_{\beta} \left\{ \sum_{t=1}^T \left[ y_t \log \sigma(\beta^T x_t) + (1 - y_t) \log (1 - \sigma(\beta^T x_t)) \right] + \lambda \sum_{i=1}^N |\beta_i| \right\} \quad (3)$$

241 Here,  $\lambda$  represents the strength of the L1 regularization term (with larger values  
 242 indicating stronger regularization), whereas  $\sigma$  represents the logistic model:

$$243 \quad \sigma(x) = \frac{1}{1 + e^{-x}} \quad (4)$$

244 For each voxel  $i$  in the brain, the resulting set of coefficient estimates  $\beta$ , indicates the  
 245 contribution of the activity of this voxel to the decoding decision  $\sigma(x_t)$  of the logistic

model for a whole-brain fMRI sample  $x_t$  at time point  $t$ . Over the recent years, the whole-brain lasso, as well as closely related decoding approaches (e.g., Gramfort et al., 2013; McIntosh and Lobaugh, 2004; Ryali et al., 2010), have found widespread application throughout the neuroscience literature (e.g., Chang et al., 2015, Wager et al., 2013).

## 2.5 DeepLight framework

### 2.5.1 Deep learning model

The DL model underlying DeepLight consists of three distinct computational modules, namely a feature extractor, an LSTM, and an output unit (for an overview, see Fig. 1). First, DeepLight separates each fMRI volume into a sequence of axial brain slices. These slices are then processed by a convolutional feature extractor (LeCun et al., 1995), resulting in a sequence of higher-level, and lower-dimensional, slice representations. These higher-level slice representations are fed to an LSTM (Hochreiter and Schmidhuber, 1997), integrating the spatial dependencies of the observed brain activity within and across axial brain slices. Lastly, the output unit makes a decoding decision, by projecting the output of the LSTM into a lower-dimensional space, spanning the cognitive states in the data. Here, a probability for each cognitive state is estimated, indicating whether the input fMRI volume belongs to each of these states. This combination of convolutional and recurrent DL elements is inspired by previous research, showing that it is generally well-suited to learn the spatial dependency structure of long sequences of input data (Donahue et al., 2015; Marban et al., 2019; McLaughlin et al., 2016). Importantly, the DeepLight approach is not dependent on any specific architecture of each of these three modules. The DL model architecture described in the following is exemplary and derived from previous work (Marban et al., 2019). Further research is needed to explore the effect of specific module architectures on the performance of DeepLight.

The feature extractor used here was composed of a sequence of eight convolution layers (LeCun et al., 1995). A convolution layer consists of a set of kernels (or filters)  $w$  that each learn local features of the input image  $a$ . These local features are then convolved over the input, resulting in an activation map  $h$ , indicating whether a feature is present at each given location of the input:

$$h_{i,j} = g \left( \sum_{k=1}^m \sum_{l=1}^m (w_{k,l} a_{i+k+1,j+l-1}) + b \right) \quad (5)$$

Here,  $b$  represents the bias of the kernel, while  $g$  represents the activation function.  $k$  and  $l$  represent the row and column index of the kernel matrix, whereas  $i$  and  $j$  represent the row and column index of the activation map.

Generally, lower-level convolution kernels (that are close to the input data) have small receptive fields and are only sensitive to local features of small patches of the input data (e.g., contrasts and orientations). Higher-level convolution kernels, on the other hand, act upon a higher-level representation of the input data, which has already been transformed by a sequence of preceding lower-level convolution kernels. Higher-level



286 kernels thereby integrate the information provided by lower-level convolution kernels,  
 287 allowing them to identify larger and more complex patterns in the data. We specified the  
 288 sequence of convolution layers as follows (see Fig. 1): conv3-16, conv3-16, conv3-16,  
 289 conv3-16, conv3-32, conv3-32, conv3-32, conv3-32 (notation: conv(kernel size) -  
 290 (number of kernels)). All convolution kernels were activated through a rectified linear  
 291 unit function:

$$292 \quad g(z) = \max(0, z) \quad (6)$$

293 Importantly, all kernels of the even-numbered convolution layers were moved over the  
 294 input fMRI slice with a stride size of one voxel and all kernels of odd-numbered layers  
 295 with a stride size of two voxels. The stride size determines the dimensionality of the  
 296 outputted slice representation. An increasing stride indicates more distance between the  
 297 application of the convolution kernels to the input data. Thereby, reducing the  
 298 dimensionality of the output representation at the cost of a decreasing sensitivity to  
 299 differences in the activity patterns of neighbouring voxels. Yet, the activity patterns of  
 300 neighbouring voxels are known to be highly correlated, leading to an overall low risk of  
 301 information loss through a reasonable increase in stride size. To avoid any further loss  
 302 of dimensionality between the convolution layers, we applied zero-padding. Thereby,  
 303 adding zeros to the borders of the inputs to each convolution layer so that the outputs of  
 304 the convolution layers have the same dimensionality as their inputs, if a stride of 1 voxel  
 305 is applied, and only decrease in size, when a larger stride is used. The sequence of eight  
 306 convolution layers thereby resulted in a 960-dimensional representation of each volume  
 307 slice.

308 To integrate the information provided by the resulting sequence of slice representations  
 309 into a higher-level representation of the observed whole-brain activity, DeepLight  
 310 applies a bi-directional LSTM (Hochreiter and Schmidhuber, 1997), containing two  
 311 independent LSTM units. Each of the two LSTM units iterates through the entire  
 312 sequence of input slices, but in reverse order (one from bottom-to-top and the other  
 313 from top-to-bottom). An LSTM unit contains a hidden cell state  $C$ , storing information  
 314 over an input sequence of length  $S$  with elements  $a_s$  and outputs a vector  $h_s$  for each  
 315 input at sequence step  $s$ . The unit has the ability to add and remove information from  $C$   
 316 through a series of gates. In a first step, the LSTM unit decides what information from  
 317 the cell state  $C$  is removed. This is done by a fully-connected logistic layer, the forget  
 318 gate  $f$ :

$$319 \quad f_t = \sigma(W_f a_s + U_f h_{s-1} + b_f) \quad (7)$$

320 Here,  $\sigma$  indicates the logistic function (see eq. 4),  $[W, U]$  the gate's weight matrices and  
 321  $b$  the gate's bias. The forget gate outputs a number between 0 and 1 for each entry in the  
 322 cell state  $C$  at the previous sequence step  $s - 1$ . Next, the LSTM unit decides what  
 323 information is going to be stored in the cell state. This operation contains two elements:  
 324 the input gate  $i$ , which decides which values of  $C_s$  will be updated, and a tanh layer,  
 325 which creates a new vector of candidate values  $C'_s$ :

$$i_s = \sigma(W_i a_s + U_i h_{s-1} + b_i) \quad (8)$$

$$C'_s = \tanh(W_c a_s + U_c h_{s-1} + b_c) \quad (9)$$

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (10)$$

Subsequently, the old cell state  $C_{s-1}$  is updated into the new cell state  $C_s$ :

$$C_s = f_s \cdot C_{s-1} + i_s \cdot C'_s \quad (11)$$

Lastly, the LSTM computes its output  $h_s$ . Here, the output gate  $o_s$ , decides what part of  $C_s$  will be outputted. Subsequently,  $C_s$  is multiplied by another  $\tanh$  layer to make sure that  $h_s$  is scaled between -1 and 1:

$$o_s = \sigma(W_o a_s + U_o h_{s-1} + b_o) \quad (12)$$

$$h_s = o_s \cdot \tanh(C_s) \quad (13)$$

Each of the two LSTM units in our DL model contained 40 output neurons. To make a decoding decision, both LSTM units pass their output for the last sequence element to a fully-connected softmax output layer. The output unit contains one neuron per cognitive state in the data and assigns a probability to each of the  $K$  (here,  $K=4$ ) states, indicating the probability that the current fMRI sample belongs to this state:

$$\sigma = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}, \text{ with } j=1, \dots, K \quad (14)$$

## 2.5.2 Layer-Wise Relevance Propagation in the DeepLight framework

To relate the decoded cognitive state and brain activity, DeepLight utilizes the Layer-Wise Relevance Propagation (LRP; Bach et al., 2015, Lapuschkin et al., 2019; Montavon et al., 2017) method. The goal of LRP is to identify the contribution of a single dimension  $d$  of an input  $a$  (with dimensionality  $D$ ) to the prediction  $f(a)$  that is made by a linear or non-linear classifier  $f$ . We denote the contribution of a single dimension as its relevance  $R_d$ . One way of decomposing the prediction  $f(a)$  is by the sum of the relevance values of each dimension of the input:

$$f(a) \approx \sum_{d=1}^D R_d \quad (15)$$

Qualitatively, any  $R_d < 0$  can be interpreted as evidence against the presence of a classification target, while  $R_d > 0$  denotes evidence for the presence of the target. Importantly, LRP assumes that  $f(a) > 0$  indicates evidence for the presence of a target.

Let's assume the relevance  $R_j^{(l)}$  of a neuron  $j$  at network layer  $l$  for the prediction  $f(a)$  is known. We would like to decompose this relevance into the messages  $R_{i \leftarrow j}^{(l-1, l)}$  that are sent to those neurons  $i$  in layer  $l-1$  which provide the inputs to neuron  $j$ :

$$R_j^{(l)} = \sum_{i \in (l)} R_{i \leftarrow j}^{(l-1, l)} \quad (16)$$

While the relevance of the output neuron at the last layer  $L$  is defined as  $R_d^{(L)} = f(a)$ , the dimension-wise relevance scores on the input neurons are given by  $R_d^{(1)}$ . For all weighted connections of the DL model in between (see eqs. 5, 7, 8, 9 and 12), DeepLight defines the messages  $R_{i \leftarrow j}^{(l-1, l)}$  as follows:

$$R_{i \leftarrow j}^{(l-1, l)} = \frac{z_{ij}}{z_j + \epsilon \cdot \text{sign}(z_j)} R_j^{(l)} \quad (17)$$

Here,  $z_{ij} = a_i^{(l-1)} w_{ij}^{(l-1, l)}$  ( $w$  indicating the coefficient weight and  $a$  the input) and  $z_j = \sum_i z_{ij}$ , while  $\epsilon$  represents a stabilizer term that is necessary to avoid numerical degenerations when  $z_j$  is close to 0 (we set  $\epsilon = 0.001$ ).

Importantly, the LSTM also applies another type of connection, which we refer to as multiplicative connection (see eqs. 11 and 13). Let  $z_j$  be an upper-layer neuron whose value in the forward pass is computed by multiplying two lower-layer neuron values  $z_g$  and  $z_s$  such that  $z_j = z_g \cdot z_s$ . These multiplicative connections occur when we multiply the outputs of a *gate* neuron, whose values range between 0 and 1, with an instance of the hidden cell state, which we will call *source* neuron. For these types of connections, we set the relevances of the gate neuron  $R_g^{(l-1)} = 0$  and the relevances of the source neuron  $R_s^{(l-1)} = R_j^{(l)}$ , where  $R_j^{(l)}$  denotes the relevances of the upper layer neuron  $z_j$  (as proposed in Arras et al., 2017). The reasoning behind this rule is that the gate neuron already decides in the forward pass how much of the information contained in the source neuron should be retained to make the classification. Even if this seems to ignore the values of the neurons  $z_g$  and  $z_s$  for the redistribution of relevance, these are actually taken into account when computing the value  $R_j^{(l)}$  from the relevances of the next upper-layer neurons to which  $z_j$  is connected by the weighted connections. We refer the reader to Samek et al. (2018) and Montavon et al. (2018) for more information about explanation methods.

In the context of this work, we decomposed the predictions of DeepLight for the actual cognitive state underlying each fMRI sample, as we were solely interested in understanding what DeepLight used as evidence in favor of the presence of this state. We also restricted the LRP analysis to those brain samples that the DL model classified correctly, because we can only assume that the DL model has learned a meaningful mapping between brain data and cognitive state, if it is able to accurately decode the cognitive state.

### 2.5.3 DeepLight training

We iteratively trained DeepLight through backpropagation (Rumelhart et al., 1986) over 60 epochs by the use of the ADAM optimization algorithm as implemented in tensorflow 1.4 (Abadi et al., 2016). To prevent overfitting, we applied dropout regularization to all network layers (Srivastava et al., 2014), global gradient norm clipping (with a clipping threshold of 5; Pascanu et al., 2013), as well as an early stopping of the training (for an overview of training statistics, see Supplementary Fig.

396 S2). During the training, we set the dropout probability to 50% for all network layers,  
397 except for the first four convolution layers, where we reduced the dropout probability to  
398 30% for the first two layers and 40% for the third and fourth layer. Each training epoch  
399 was defined as a complete iteration over all samples in the training dataset (see Section  
400 2.2). We used a learning rate of 0.0001 and a batch size of 32. All network weights were  
401 initialized by the use of a normal-distributed random initialization scheme (Glorot and  
402 Bengio, 2010). The DL model was written in tensorflow 1.4 (Abadi et al., 2016) and the  
403 interprettensor library (<https://github.com/VigneshSrinivasan10/interprettensor>).

#### 404 **2.5.4 DeepLight brain maps**

405 To generate a set of subject-level brain maps with DeepLight, we first decomposed the  
406 decoding decisions of DeepLight for each correctly classified fMRI sample of a subject  
407 with the LRP method (see Section 2.5.2). Importantly, we restricted the LRP analysis to  
408 those fMRI samples that were collected 5 - 15s after the onset of the experiment block,  
409 as we expect the HRF (Lindquist et al., 2009) to be strongest within this time period. To  
410 then aggregate the resulting set of relevance maps for each decomposed fMRI sample  
411 within each cognitive state, we smoothed each relevance map with a 3mm FWHM  
412 Gaussian kernel and averaged all relevance volumes belonging to a cognitive state,  
413 resulting in one brain map per subject and cognitive state. Group-level brain maps were  
414 then obtained, by averaging these subject-level brain maps for all subjects in the held-  
415 out test dataset within each cognitive state, resulting in one group-level brain map per  
416 cognitive state.

### 417 **3. Results**

#### 418 **3.1 DeepLight accurately decodes cognitive states from fMRI data**

419 A key prerequisite for the DeepLight analysis (as well as all other decoding analyses) is  
420 that it achieves reasonable performance in the decoding task at hand. Only then we can  
421 assume that it has learned a meaningful mapping from the fMRI data to the cognitive  
422 states and interpret the resulting brain maps as informative about these states.

423 Overall, DeepLight accurately decoded the cognitive states underlying 68.3% of the  
424 fMRI samples in the held-out test dataset (62.36%, 69.87%, 75.97%, 65.09% for body,  
425 face, place and tool respectively; Fig. 2A). It generally performed best at discriminating  
426 the body and place (5.1% confusion in the held-out data), face and tool (7.8% confusion  
427 in the held-out data), body and tool (9.8% confusion in the held-out data) and face and  
428 place (10.4% confusion in the held-out data) stimuli from the fMRI data, while it did  
429 not perform as well in discriminating place and tool and body and face stimuli (15%  
430 confusion in the held-out data respectively).

431 Note that DeepLight's performance in decoding the four cognitive states from the fMRI  
432 data varied over the course of an experiment block (Fig. 2B). DeepLight performed best  
433 in the middle and later stages of the experiment block, where the average decoding  
434 accuracy reaches 80%. This finding is generally in line with the temporal evolution of  
435 the hemodynamic response function (HRF; Lindquist et al., 2009) measured by the

436 fMRI (the HRF is known to be strongest 5-10 seconds after to the onset of the  
437 underlying neuronal activity).

438 To further evaluate DeepLight's performance in decoding the cognitive states from the  
439 fMRI data, we compared its performance in decoding these states to the searchlight  
440 analysis and whole-brain lasso. For simplicity, we sub-divided this comparison into a  
441 separate analysis on the group- and subject-level.

### 442 3.1.1 Group-level

443 For the group-level comparison, we trained the searchlight analysis and whole-brain  
444 lasso on the data of all 70 subjects contained in the training dataset (for details on the  
445 fitting procedures, see Supplementary Information Section 1). Subsequently, we  
446 evaluated their performance in decoding the cognitive states in the full held-out test  
447 data.

448 DeepLight clearly outperformed the other approaches in decoding the cognitive states.  
449 While the searchlight analysis achieved an average decoding accuracy of 60% (Fig. 2C)  
450 and the whole-brain lasso an average decoding accuracy of 47.97% (Fig. 2D),  
451 DeepLight improved upon these performances by 8.3% ( $t(29)=5.80$ ,  $p<0.0001$ ) and  
452 20.33% ( $t(29)=13.39$ ,  $p<0.0001$ ) respectively.

453 All three decoding approaches generally performed best at discriminating face and place  
454 stimuli from the fMRI data (Fig. 2A, C-D). Similar to DeepLight, the searchlight  
455 analysis and whole-brain lasso also performed well at discriminating body and place  
456 stimuli (3.3% and 12.2% confusion for the searchlight analysis and whole-brain lasso  
457 respectively, Fig. 2C-D), while they also had more difficulties discriminating body and  
458 face stimuli from the fMRI data (25% and 20.2% confusion for the searchlight analysis  
459 and whole-brain lasso respectively, Fig. 2C-D).

460 A key premise of DL methods, when compared to more traditional decoding  
461 approaches, is that their decoding performance improves better with growing datasets.  
462 To test this, we repeatedly trained all three decoding approaches on a subset of the  
463 training dataset (including the data of 5, 10, 15, 20, 25, 30, 35, 40, 50, 60 and 70  
464 subjects), and validated their performance at each iteration on the full held-out test data  
465 (Fig. 2E). Overall, the decoding performance of DeepLight increased by 0.27%  
466 ( $t(10)=10.9$ ,  $p<0.0001$ ) per additional subject in the training dataset, whereas the  
467 performance of the whole-brain lasso increased by 0.03% ( $t(10)=3.02$ ,  $p=0.015$ ) and the  
468 performance of the searchlight analysis only marginally increased by 0.04%  
469 ( $t(10)=2.08$ ,  $p=0.067$ ). Nevertheless, the searchlight analysis outperformed DeepLight  
470 in decoding the cognitive states from the data when only little training data were  
471 available (here, 10 or less subjects ( $t(29)=-4.39$ ,  $p<0.0001$ ). The decoding advantage of  
472 DeepLight, on the other hand, came to light when the data of 50 or more subjects were  
473 available in the training dataset ( $t(29)=3.82$ ,  $p=0.0006$ ). DeepLight consistently  
474 outperformed the whole-brain lasso, when it was trained on the data of at least 10  
475 subjects ( $t(29)=5.32$ ,  $p=0.0045$ ).

### 3.1.2 Subject-level

For the subject-level comparison, we first trained both, the searchlight analysis and whole-brain lasso on the fMRI data of the first experiment run of a subject from the held-out test dataset (for an overview of the training procedures, see Supplementary Information Section 1). We then used the data of the second experiment run of the same subject to evaluate their decoding performance (by predicting the cognitive states underlying each fMRI sample of the second experiment run). Importantly, we also decoded the same fMRI samples with DeepLight. Note that DeepLight, in comparison to the other approaches, did not see any data of the subject during the training, as it was solely trained on the data of the 70 subjects in the training dataset (see Section 2.1).

DeepLight clearly outperformed the other decoding approaches, by decoding the cognitive states more accurately for 28 out of 30 subjects, when compared to the searchlight analysis (while the searchlight analysis achieved an average decoding accuracy of 47.2% across subjects, DeepLight improved upon this performance by 22.4%, with an average decoding accuracy of 69.3%,  $t(29)=11.28$ ,  $p<0.0001$ ; Fig. 3A), and for 29 out of 30 subjects, when compared to the whole-brain lasso (while the whole-brain lasso achieved an average decoding accuracy of 37% across subjects, DeepLight improved upon this performance by 32%;  $t(29)=15.74$ ,  $p<0.0001$ ; Fig. 3B).

To further ascertain that the observed differences in decoding performance between the searchlight and DeepLight did not result from the linearity contained in the Support Vector Machine (SVM; Cortes and Vapnik, 1995) of the the searchlight analysis, we replicated our subject-level searchlight analysis, by the use of a non-linear radial basis function kernel (RBF; Cortes and Vapnik, 1995, Müller et al., 2001, Schölkopf and Smola, 2002) SVM (Supplementary Fig. S3). However, the decoding accuracies achieved by the RBF-kernel SVM were not meaningfully different from those of the linear-kernel SVM ( $t(29)=-1.75$ ,  $p=0.09$ ).

Lastly, we also compared the subject-level decoding performance of the whole-brain lasso to that of a recently proposed extension of this approach (TV-L1, for methodological details see Gramfort et al., 2013). The TV-L1 approach combines the Least Absolute Shrinkage Regularization (L1; see eq. 3) of the whole-brain lasso with an additional Total-Variation (TV) penalty (Michel et al., 2011), to better account for the spatial dependency structure of fMRI data. Yet, we found that the whole-brain lasso performed better at decoding the cognitive states from the subject-level fMRI data than TV-L1 ( $t(29)=3.79$ ,  $p=0.0007$ ; see Supplementary Fig. S4).

### 3.2 DeepLight identifies physiologically appropriate associations between cognitive states and brain activity

Our previous analyses have shown that DeepLight has learned a meaningful mapping between the fMRI data and cognitive states, by accurately decoding these states from the data. Next, we therefore tested DeepLight's ability to identify the brain areas associated with the cognitive states, by decomposing its decoding decisions with the LRP method (see Section 2.5). Subsequently, we compared the resulting brain maps of DeepLight to those of the GLM, searchlight analysis and whole-brain lasso. Again, we

sub-divided this comparison into a separate analysis on the group- and subject-level. Note that due to the diverse statistical nature of the three baseline approaches, the values of their brain maps are on different scales and have different statistical interpretations (for methodological details, see Section 2.4). Further, all depicted brain maps in Fig. 4-6 are projected onto the inflated cortical surface of the FsAverage5 surface template (Fischl, 2012) for better visibility.

To evaluate the quality of the brain maps resulting from each analysis approach, we performed a meta-analysis of the four cognitive states with NeuroSynth (for details on NeuroSynth, see Supplementary Information Section 2 and Yarkoni et al., 2011). NeuroSynth provides a database of mappings between cognitive states and brain activity, based on the empirical neuroscience literature. Particularly, the resulting brain maps used here indicate whether the probability that an article reports a specific brain activation is different, when it includes a specific term (e.g., "face") compared to when it does not. With this meta-analysis, we defined a set of *regions-of-interest* (ROIs) for each cognitive state (as defined by the terms "body", "face", "place", and "tools"), in which we would expect the various analysis approaches to identify a positive association between the cognitive state and brain activity (for an overview, see Fig. 4A). These ROIs were defined as follows: the upper parts of the middle and inferior temporal gyrus, the postcentral gyrus, as well as the right fusiform gyrus for the body state, the fusiform gyrus (also known as the fusiform face area FFA; Haxby et al., 2001, Heekeren et al., 2004) and amygdala for the face state, the parahippocampal gyrus (or parahippocampal place area PPA; Haxby et al., 2001, Heekeren et al., 2004) for the place state and the upper left middle and inferior temporal gyrus as well as the left postcentral gyrus for the tool state.

To ensure comparability with the results of the meta-analysis, we restricted all analyses of brain maps to the estimated positive associations between brain activity and cognitive states (i.e., positive relevance values as well as positive GLM and whole-brain lasso coefficients, see Section 2.4 and Supplementary Information Section 1). A negative Z-value in the meta-analysis indicates a lower probability that an article reports a specific brain activation when it includes a specific term, compared to when it does not include the term. A negative value in the meta-analysis is therefore conceptually different to negative values in the brain maps of our analyses (e.g., negative relevance values or negative whole-brain lasso coefficients). These can generally be interpreted as evidence against the presence of a cognitive state, given the specific set of cognitive states in our dataset (e.g., a negative relevance indicates evidence for the presence of any of the other cognitive states considered).

### 3.2.1 Group-level

To determine the voxels that each analysis approach associated with a cognitive state, we defined a threshold for the values of each group-level brain map, indicating those voxels that are associated most strongly with the cognitive state. For the GLM analysis, we thresholded all P-values at an expected false discovery rate (Benjamini & Hochberg, 1995; Genovese, Lazar & Nichols, 2002) of 0.1 (Fig. 4B). Similarly, for all decoding analyses, we thresholded each brain map at the 90th percentile of its values (Fig. 4C-E). For the whole-brain lasso and DeepLight, the remaining 10 percent of values indicate

562 those brain regions whose activity these approaches generally weight most in their  
563 decoding decisions. For the searchlight analysis, the remaining 10 percent of values  
564 indicate those brain regions in which the searchlight analysis achieved the highest  
565 decoding accuracy.

566 All analysis approaches correctly associated activity in the upper parts of the middle and  
567 inferior temporal gyrus with body stimuli. The GLM, whole-brain lasso and DeepLight  
568 also correctly associated activity in the right fusiform gyrus with body stimuli. Only  
569 DeepLight correctly associated activity in the postcentral gyrus with these stimuli. The  
570 GLM, whole-brain lasso and DeepLight further all correctly associated activity in the  
571 right FFA with face stimuli. None of the approaches, however, associated activity in the  
572 left FFA with face stimuli. Interestingly, the searchlight analysis did not associate the  
573 FFA with face stimuli at all. All analysis approaches also correctly associated activity in  
574 the PPA with place stimuli. Lastly, for tool stimuli, the GLM and whole-brain lasso  
575 correctly associated activity in the left inferior temporal sulcus with stimuli of this class.  
576 The searchlight analysis and whole-brain lasso only did so marginally. None of the  
577 approaches associated activity in the left postcentral gyrus with tool stimuli.

578 Overall, DeepLight's group-level brain maps accurately associated each of the ROIs  
579 with their respective cognitive states. Interestingly, DeepLight also associated a set of  
580 additional brain regions with the face and tool stimulus classes that were not identified  
581 by the other analysis approaches (see Fig. 4E). For face stimuli, these regions are the  
582 orbitofrontal cortex and temporal pole. While the temporal pole has been shown to be  
583 involved in the ability of an individual to infer the desires, intentions and beliefs of  
584 others (*theory-of-mind*; for a detailed review, see Olson et al., 2007), the orbitofrontal  
585 cortex has been associated with the processing of emotions in the faces of others (for a  
586 detailed review, see Adolphs, 2002). For tool stimuli, DeepLight additionally utilized  
587 the activity of the temporoparietal junction (TPJ) to decode these stimuli. The TPJ has  
588 been shown to be associated with the ability of an individual to discriminate self-  
589 produced actions and the actions produced by others and is generally regarded of as a  
590 central hub for the integration of body-related information (for a detailed review, see  
591 Decety and Grèzes, 2006). Although it is not clear why only DeepLight associated these  
592 brain regions with the face and tool stimulus classes, their assumed functional roles do  
593 not contradict this association.

### 594 3.2.2 Subject-level

595 The goal of the subject-level analysis was to test the ability of each analysis approach to  
596 identify the physiologically appropriate associations between brain activity and  
597 cognitive state on the level of each individual.

598 To quantify the similarity between the subject-level brain maps and the results of the  
599 meta-analysis, we defined a similarity measure. Given a target brain map (e.g., the  
600 results of our meta-analysis), this measure tests for each voxel in the brain whether a  
601 source brain map (e.g., the results of our subject-level analyses) correctly associates this  
602 voxel's activity with the cognitive state (true positive), falsely associates the voxel's  
603 activity with the cognitive state (false positives) or falsely does not associate the voxel's  
604 activity with the cognitive state (false negatives). Particularly, we derived this measure



605 from the well-known F1-score in machine learning (see Supplementary Information  
606 Section 3 as well as Goutte and Gaussier, 2005). The benefit of the F1-score, when  
607 compared to simply computing the ratio of correctly classified voxels in the brain, is  
608 that it specifically considers the brain map's precision and recall and is thereby robust to  
609 the overall size of the ROIs in the target brain map. Here, precision describes the  
610 fraction of true positives from the total number of voxels that are associated with a  
611 cognitive state in the source brain map. Recall, on the other hand, describes the fraction  
612 of true positives from the overall number of voxels that are associated with a cognitive  
613 state in the target brain map. Generally, an F1-score of 1 indicates that the brain map  
614 has both, perfect precision and recall with respect to the target, whereas the F1-score is  
615 worst at 0.

616 To obtain an F1-score for each subject-level brain map (for details on the estimation of  
617 subject-level brain maps with the three baseline analysis approaches, see Supplementary  
618 Information Section 1), we again thresholded each individual brain map. For the GLM,  
619 we defined all voxels with a P-value greater than 0.005 (uncorrected) as not associated  
620 with the cognitive state and all others as associated with the cognitive state. For the  
621 searchlight analysis, whole-brain lasso and DeepLight, we defined all voxels with a  
622 value below the 90th percentile of the values within the brain map as not associated with  
623 the cognitive state and all others as associated with the cognitive state.

624 Overall, DeepLight's subject-level brain maps had meaningfully larger F1-scores for the  
625 body, face and place stimulus classes, when compared to those of the GLM  
626 ( $t(29)=10.46$ ,  $p<0.0001$  for body stimuli, Supplementary Fig. S5A;  $t(29)=13.04$ ,  
627  $p<0.0001$  for face stimuli, Supplementary Fig. S5D;  $t(29)=9.26$ ,  $p<0.0001$  for place  
628 stimuli, Supplementary Fig. S5G), searchlight analysis ( $t(29)=13.26$ ,  $p<0.0001$  for body  
629 stimuli, Supplementary Fig. S5B;  $t(29)=8.57$ ,  $p<0.0001$  for face stimuli, Supplementary  
630 Fig. S5E;  $t(29)=4.25$ ,  $p=0.0002$ , for place stimuli, Supplementary Fig. S5H), and  
631 whole-brain lasso ( $t(29)=20.93$ ,  $p<0.0001$  for body stimuli, Supplementary Fig. S5C;  
632  $t(29)=48.32$ ,  $p<0.0001$  for face stimuli, Supplementary Fig. S5F;  $t(29)=22.43$ ,  
633  $p<0.0001$ , for place stimuli, Supplementary Fig. S5I). For tool stimuli, the GLM and  
634 searchlight generally achieved higher subject-level F1-scores than DeepLight ( $t(29)=-$   
635  $8.19$ ,  $p<0.0001$ , Supplementary Fig. S5J;  $t(29)=-4.39$ ,  $p=0.0001$ , Supplementary Fig.  
636 S5K for the GLM and searchlight respectively), whereas DeepLight outperformed the  
637 whole-brain lasso analysis ( $t(29)=18.31$ ,  $p<0.0001$ , Supplementary Fig. S5L).

638 To ascertain that the results of this comparison were not dependent on the thresholds  
639 that we chose, we replicated the comparison for each combination of the 85th, 90th and  
640 95th percentile threshold for the brain maps of the searchlight analysis, whole-brain  
641 lasso and DeepLight, as well as a P-threshold of 0.05, 0.005, 0.0005 and 0.00005 for the  
642 brain maps of the GLM. Within all combinations of percentile values and P-thresholds,  
643 the presented results of the F1-comparison were generally stable (see Supplementary  
644 Table S3-6).

### 3.3 DeepLight accurately identifies physiologically appropriate associations between cognitive states and brain activity on multiple levels of data granularity

DeepLight's ability to correctly identify the physiological appropriate associations between cognitive states and brain activity is exemplified in Figure 5. Here, the distribution of relevance values for the four cognitive states is visualized on three different levels of data granularity of an exemplar subject (namely, the subject with the highest decoding accuracy in Fig. 3A-B): First, on the level of the overall distribution of relevance values of each cognitive state of this subject (Fig. 5A; incorporating an average of 47 TRs per cognitive state), then on the level of the first experiment block of each cognitive state in the first experiment run (Fig. 5B; incorporating an average of 12 TRs per cognitive state) and lastly on the level of a single brain sample of each cognitive state (Fig. 5C; incorporating a single TR per cognitive state).

On all three levels, DeepLight utilized the activity of a similar set of brain regions to identify each of the four cognitive states. Importantly, these regions largely overlap with those identified by the DeepLight group-level analysis (Fig. 4E) as well as the results of the meta-analysis (Fig. 4A).

### 3.4 DeepLight's relevance patterns resemble temporo-spatial variability of brain activity over sequences of single fMRI samples

To further probe DeepLight's ability to analyze single time points, we next studied the distribution of relevance values over the course of a single experiment block (Fig. 6). In particular, we plotted this distribution as a function of the fMRI sampling-time over all subjects for the first experiment block of the face and place stimulus classes in the second experiment run. We restricted this analysis to the face and place stimulus classes, as the neural networks involved in processing face and place stimuli, respectively, have been widely characterized (see, for example Haxby et al., 2001 as well as Heekeren et al., 2004). For a more detailed overview, we also created two videos for the two experiment blocks depicted in Figure 6 (Supplementary Videos 1 and 2). These videos display the temporal evolution of relevance values for each fMRI sample in the original fMRI sampling time of the face (Supplementary Video 1) and place (Supplementary Video 2) experiment blocks.

In the beginning of the experiment block, DeepLight was generally uncertain which cognitive state the observed brain samples belonged to, as it assigned similar probabilities to each of the cognitive states considered (Fig. 6A-B). As time progressed, however, DeepLight's certainty increased and it correctly identified the cognitive state underlying the fMRI samples. At the same time, it started assigning more relevance to the target ROIs of the face and place stimulus classes (Fig. 6C-F), as indicated by the increasing F1-scores resulting from a comparison of the brain maps at each sampling time point with the results of the meta-analysis (Fig. 6G-H; all brain maps were again thresholded at the 90th percentile for this comparison). Interestingly, the relevances started peaking in the target ROIs 5s after the onset of the experiment block. The

686 temporal evolution of the relevances thereby mimics the hemodynamic response  
687 measured by the fMRI (Lindquist et al., 2009).

688 To further evaluate the results of this analysis, we replicated it by the use of the whole-  
689 brain lasso group-level decoding model (see Section 2.4 and Supplementary Information  
690 Section 1). In particular, we multiplied the fMRI samples of all test subjects collected at  
691 each sampling time point with the coefficient estimates of the whole-brain lasso group-  
692 level model. Subsequently, we averaged the resulting weighted fMRI samples within  
693 each sampling time point depicted in Fig. 6G-H and computed an F1-score for a  
694 comparison of the resulting average brain maps with the results of the meta-analysis (as  
695 described in section 3.2.2). Interestingly, we found that the F1-scores of the whole-brain  
696 lasso analysis varied much less over the sequence of fMRI samples and were throughout  
697 lower than those of DeepLight. Thereby, indicating that the brain maps of the whole-  
698 brain lasso analysis exhibit comparably little variability over the course of an  
699 experiment block with respect to the target ROIs defined for the face and place stimulus  
700 classes.

## 701 **4. Discussion**

702 Neuroimaging data have a complex temporo-spatial dependency structure that renders  
703 modeling and decoding of experimental data a challenging endeavor. With DeepLight,  
704 we propose a new data-driven framework for the analysis and interpretation of whole-  
705 brain neuroimaging data that scales well to large datasets and is mathematically non-  
706 linear, while still maintaining interpretability of the data. To decode a cognitive state,  
707 DeepLight separates a whole-brain fMRI volume into its axial slices and processes the  
708 resulting sequence of brain slices by the use of a convolutional feature extractor and  
709 LSTM. Thereby, accounting for the spatially distributed patterns of whole-brain brain  
710 activity within and across axial slices. Subsequently, DeepLight relates cognitive state  
711 and brain activity, by decomposing its decoding decisions into the contributions of the  
712 single input voxels to these decisions with the LRP method. Thus, DeepLight is able to  
713 study the associations between brain activity and cognitive state on multiple levels of  
714 data granularity, from the level of the group down to the level of single subjects, trials  
715 and time points.

716 To demonstrate the versatility of DeepLight, we have applied it to an openly available  
717 fMRI dataset of 100 subjects viewing images of body parts, faces, places and tools.  
718 With these data, we have shown that the DeepLight 1) decodes the underlying cognitive  
719 states more accurately from the fMRI data than conventional means of uni- and  
720 multivariate brain decoding, 2) improves its decoding performance better with growing  
721 datasets, 3) accurately identifies the physiologically appropriate associations between  
722 cognitive states and brain activity, 4) can study these associations on multiple levels of  
723 data granularity, from the level of the group down to the level of single subjects, trials  
724 and time points and 5) can capture the temporo-spatial variability of brain activity over  
725 sequences of single fMRI samples.

## 726 4.1 Transferring DeepLight to other fMRI datasets

727 The DeepLight architecture used here is exemplary. Future research is needed to  
728 evaluate how the specific architectural choices for its three sub-modules (the  
729 convolutional feature extractor, LSTM unit and softmax output layer; see Section 2.5)  
730 will effect its performance. In the following, we will briefly outline how the proposed  
731 architecture can be transferred to the analysis of other fMRI datasets with different  
732 spatial resolution and decoding targets. Importantly, online minimal changes are  
733 necessary in order to adapt DeepLight’s architecture for the analysis of such fMRI  
734 datasets.

735 DeepLight first processes an fMRI volume within each axial slice, by computing a  
736 higher-level, and lower-dimensional, representation of the slices with the convolutional  
737 feature extractor. Here, the spatial sensitivity of DeepLight to the fine-grained activity  
738 differences of neighboring voxels within each slice is determined by the stride size  
739 applied by the convolution layers. The stride size indicates the distance between the  
740 application of the convolution kernels to the axial slices of the fMRI volume (see eq. 5).  
741 Generally, a larger stride decreases DeepLight’s sensitivity for fine-grained differences  
742 in the activity of neighboring voxels, as it increases the distance between the  
743 applications of the convolution kernels to the input slice. Reversely, a smaller stride size  
744 increases DeepLight’s sensitivity for the fine-grained activity differences of neighboring  
745 voxels, as it decreases the distance between the applications of the convolution kernels.  
746 For example, when analyzing fMRI volumes that have a lower spatial resolution than  
747 the ones used here, containing fewer voxels per axial slice (and thereby less information  
748 about the distribution of brain activity within each slice), we would recommend to  
749 decrease the stride size for more of DeepLight’s convolution layers, in order to best  
750 leverage the information contained in these voxels.

751 After the application of the convolutional feature extractor, DeepLight integrates the  
752 information of the resulting higher-level slice representations, by the use of a bi-  
753 directional LSTM. Here, each of the two LSTM units iterates through the entire  
754 sequence of slice representations, before forwarding its output. The proposed DeepLight  
755 architecture therefore does not require any modification in order to accommodate fMRI  
756 datasets with a different number of axial slices per volume, as it generalizes to any  
757 sequence length.

758 Further, the number of neurons in the softmax output layer is directly determined by the  
759 number of decoding targets considered in the data (one output neuron per decoding  
760 target). In the case of a continuous decoding target (for example, by predicting a  
761 subject’s score in a cognitive test), the softmax output layer can be replaced with a  
762 linear regression layer. The LRP decomposition approach (see Section 2.5.2) also  
763 applies to continuous output variables (for further details on the application of the LRP  
764 approach to continuous output variables, see Bach et al., 2015 and Montavon et al.,  
765 2017).

766 Lastly, recent exploratory empirical work has shown that even for more complex fMRI  
767 decoding analyses, encompassing up to 400 subjects and 20 distinct cognitive states (see  
768 Thomas et al., 2019), DeepLight does not require more than 64 neurons per layer. We

769 would therefore not recommend to increase the number of neurons further, as this will  
770 also lead to an overall increased risk of overfitting.

## 771 **4.2 Comparison to baseline methods**

### 772 **4.2.1 General linear model**

773 The GLM is conceptually different from the other neuroimaging analysis approaches  
774 considered in this work. It aims to identify an association between cognitive state and  
775 brain activity, by modeling (or predicting) the time series signal of a single voxel as a  
776 linear combination of a set of experiment predictors (see Section 2.4). It is thereby  
777 limited in three meaningful ways that do not apply to DeepLight: First, the time series  
778 signal of a voxel is generally very noisy. The GLM treats each voxel's signal as  
779 independent of one another, thereby, not leveraging the evidence that is shared across  
780 the time series signal of multiple voxels. Second, even though the linear combination of  
781 a set of experiment predictors might be able to explain variance in the observed fMRI  
782 data, it does not necessarily provide evidence that this exact set of predictors is encoded  
783 in the neuronal response. Generally, the same linear model (in terms of its predictions)  
784 can be constructed from many different (even random) sets of predictors (for a detailed  
785 discussion of this "feature fallacy", see Kriegeskorte and Douglas, 2018). The results of  
786 the GLM analysis thereby solely indicate that the measured neuronal response is highly  
787 structured and that this structure is preserved across individuals, whereas the labels  
788 assigned to its predictors might be arbitrary. Third, the performance of the GLM in  
789 predicting the response signal of a voxel is typically not evaluated on independent data,  
790 which leaves unanswered how well its results generalize to new data.

### 791 **4.2.2 Searchlight analysis**

792 DeepLight generally outperformed the searchlight analysis in decoding the cognitive  
793 states from the fMRI data. In small datasets (here, the data of 10 or less subjects),  
794 however, the performance of the searchlight analysis was superior. In contrast to  
795 DeepLight, the searchlight analysis decodes a cognitive state from single clusters of  
796 only few voxels. Its input data, as well as the number of parameters in its decoding  
797 model, are thereby considerably smaller, leading to an overall lower risk of overfitting.  
798 Yet, this advantage comes at the cost of additional constraints that have to be considered  
799 when choosing between both approaches. If a cognitive state is associated with the  
800 activity of a small brain region only, the searchlight analysis will generally be more  
801 sensitive to the activity of this region than DeepLight, as it has learned a decoding  
802 model that is specific to the activity of the region. If, however, the cognitive state is not  
803 identifiable by the activity of a single brain region only, but solely in conjunction with  
804 the activity of another spatially distinct brain region, the searchlight analysis will not be  
805 able to identify this association, due to its narrow spatial focus. DeepLight, on the other  
806 hand, will generally be less sensitive to the specifics of the activity of a local brain  
807 region, but perform better in identifying a cognitive state from spatially wide-spread  
808 brain activity. When choosing between both approaches, one should therefore consider  
809 whether the assumed associations between brain activity and cognitive state specifically  
810 involve the activity of a local brain region only, or whether the cognitive state is  
811 associated with the activity of spatially distinct brain regions.

### 812 4.2.3 Whole-brain lasso

813 In contrast to DeepLight, the whole-brain lasso analysis is based on a linear decoding  
814 model. It assigns a single coefficient weight to each voxel in the brain and makes a  
815 decoding decision by computing a weighted sum over the activity of an input fMRI  
816 volume. Importantly, due to the strong regularization that is applied to the coefficients  
817 during the training, many coefficients equal 0. The resulting set of coefficients thereby  
818 resembles a brain mask, defining a set of fixed brain regions whose activity the whole-  
819 brain lasso utilizes to decode a cognitive state. DeepLight, on the other hand, utilizes a  
820 hierarchical structure of non-linear transforms of the fMRI data. It projects each fMRI  
821 volume into a more abstracted, higher-level space. This abstracted (and more flexible)  
822 view enables DeepLight to better account for the variable patterns of brain activity  
823 underlying a cognitive state (within and across individuals). This ability is exemplified  
824 in Figure 6, as well as Supplementary Videos 1 - 2, where we visualize the variable  
825 patterns of brain activity that DeepLight associates with the face and place stimulus  
826 classes throughout an experiment block. The relevance patterns of DeepLight mimic the  
827 hemodynamic response and peak in the ROIs 5-10s after the onset of the experiment  
828 block. Importantly, we find that the whole-brain lasso does not exhibit such temporo-  
829 spatial variability.

### 830 4.3 Disentangling temporally distinct associations between cognitive 831 state and brain activity

832 DeepLight’s ability to identify a cognitive state through variable patterns of brain  
833 activity makes it ideally suited for the analysis of the fine-grained spatial distribution of  
834 brain activity over temporal sequences of fMRI samples. For example, Hunt and  
835 Hayden (2017) recently raised the question whether the neural networks involved in  
836 reward-based decision making can be subdivided into a set of spatially distinct and  
837 temporally discrete network components, or whether the underlying networks act in  
838 parallel, with highly recurrent activity patterns. Answering this question is difficult with  
839 conventional approaches to the analysis of neuroimaging data, such as the baseline  
840 methods included in this paper. These often learn a fixed mapping between brain  
841 activity and cognitive state, by aggregating over the information provided by a sequence  
842 of fMRI samples (e.g., by estimating a single coefficient weight for each voxel from a  
843 sequence of fMRI data). The resulting brain maps thereby only indicate whether there  
844 exist spatially distinct brain regions that are associated with a cognitive state, without  
845 providing any insight whether the activity patterns are temporally discrete. While these  
846 methods can be adapted to specifically account for the temporal differences in the  
847 activity patterns of these regions (e.g., by analyzing different time points independent of  
848 one another), these adaptations often require specific hypotheses about the studied  
849 temporal differences (e.g., by needing to specify the different time points to analyze).  
850 DeepLight, on the other hand, operates purely data-driven and is thereby able to  
851 autonomously identify an association between spatially distinct patterns of brain activity  
852 and a cognitive state at temporally discrete time points.

#### 853 **4.4 Integrative analysis of multimodal neuroimaging data**

854 DeepLight is not bound to fMRI data, but can be easily extended to other neuroimaging  
855 modalities. One such complementary modality, with a higher temporal, but lower spatial  
856 resolution, is the Electroencephalography (EEG). While a plethora of analysis  
857 approaches have been proposed for the integrative analysis of EEG and fMRI data,  
858 these often incorporate restrictive assumptions to enable the integrative statistical  
859 analysis of these two data types, with clearly distinct spatial, temporal and physiological  
860 properties (for a detailed review, see Jorge et al., 2014). DeepLight, on the other hand,  
861 represents a data-driven analysis framework. By providing both, EEG and fMRI data as  
862 separate inputs to the DL model, DeepLight could learn the fine-grained temporal  
863 structure of brain activity from the EEG data, while utilizing the fMRI data to localize  
864 the spatial brain regions underlying this activity. Recently, researchers have already  
865 demonstrated the usefulness of interpretable DL methods for the analysis of EEG data  
866 (Sturm et al., 2016).

#### 867 **4.5 Extending DeepLight**

868 Lastly, we would like to highlight several possible extensions of the DeepLight  
869 approach, resulting from its flexible and modular architecture. First, DeepLight can be  
870 extended to specifically account for the temporo-spatial distribution of brain activity  
871 over sequences of fMRI samples, by the addition of another recurrent network layer.  
872 This layer would process each of the higher-level whole-brain representations resulting  
873 from the currently proposed architecture. This extension would enable DeepLight to  
874 more specifically account for the temporal distribution of brain activity. Second,  
875 DeepLight can be extended to the integrative analysis of neuroimaging data from  
876 multiple cognitive tasks and experiments. For example, by adding one neuron to the  
877 output layer for each cognitive state from each task. This extension would enable a more  
878 thorough analysis of the differences (and similarities) between the associations of  
879 cognitive state and brain activity across multiple tasks and experiments.

#### 880 **5. Conflict of Interest**

881 The authors declare that the research was conducted in the absence of any commercial  
882 or financial relationships that could be construed as a potential conflict of interest.

#### 883 **6. Author Contributions**

884 A.W.T., H.R.H., K.R.M., and W.S. conceived of DeepLight. A.W.T., K.R.M. and W.S.  
885 planned all data analyses. A.W.T. implemented all visualizations of DeepLight and the  
886 experimental procedures and performed all formal data analyses. A.W.T. wrote all  
887 software that was used in the data analyses and that is underlying DeepLight. A.W.T.  
888 wrote the original draft of the manuscript, and H.R.H., K.R.M., and W.S. reviewed and  
889 edited the manuscript. The work was supervised by H.R.H., K.R.M., and W.S..

#### 890 **7. Funding**

891 This work was supported by the German Federal Ministry for Education and Research  
892 through the Berlin Big Data Centre (01IS14013A), the Berlin Center for Machine  
893 Learning (01IS18037I) and the TraMeExCo project (01IS18056A). Partial funding by  
894 the German Research Foundation (DFG) is acknowledged (EXC 2046/1, project-ID:  
895 390685689). KRM is also supported by the Information & Communications  
896 Technology Planning & Evaluation (IITP) grant funded by the Korea government (No.  
897 2017-0-00451).

## 898 **8. Acknowledgments**

899 This manuscript has been released as a Pre-Print at arXiv:1810.09945 (see Thomas et  
900 al., 2018).

## 901 **9. Ethics Statement**

902 The scanning protocol was approved by Washington University in St. Louis's Human  
903 Research Protection Office (HRPO), IRB# 201204036. No experimental activity  
904 involving the human subjects took place at the authors' institutions. The participants  
905 included in this study provided written informed consent and were scanned according to  
906 procedures approved by the IRB at Washington University. Only de-identified, publicly  
907 released data were used in this study.

## 908 **10. References**

- 909 Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat,  
910 S., Irving, G., Isard, M., et al. (2016). Tensorflow: A system for large-scale machine  
911 learning. In OSDI, volume 16, pages 265–283.
- 912 Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J.,  
913 Gramfort, A., Thirion, B., and Varoquaux, G. (2014). Machine learning for  
914 neuroimaging with scikit-learn. *Frontiers in Neuroinformatics*, 8:14.
- 915 Adolphs, R. (2002). Neural systems for recognizing emotion. *Current opinion in*  
916 *neurobiology*, 12(2):169–177.
- 917 Arras, L., Montavon, G., Müller, K.-R., and Samek, W. (2017). Explaining recurrent  
918 neural network predictions in sentiment analysis. In *Proceedings of the EMNLP'17*  
919 *Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media*  
920 *Analysis (WASSA)*, pages 159-168. Association for Computational Linguistics.
- 921 Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., and Samek, W.  
922 (2015). On pixel-wise explanations for non-linear classifier decisions by layer-wise  
923 relevance propagation. *PLOS ONE*, 10(7):e0130140.
- 924 Barch, D. M., Burgess, G. f., Harms, M. P., Petersen, S. E., Schlaggar, B. L., Corbetta,  
925 M., Glasser, M. F., Curtiss, S., Dixit, S., Feldt, C., et al. (2013). Function in the  
926 human connectome: task-fMRI and individual differences in behavior. *Neuroimage*,  
927 80:169–189.



- 928 Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: a practical  
929 and powerful approach to multiple testing. *Journal of the Royal statistical society:*  
930 *series B (Methodological)*, 57(1), 289-300.
- 931 Chang, L. J., Gianaros, P. J., Manuck, S. B., Krishnan, A., and Wager, T. D. (2015). A  
932 sensitive and specific neural signature for picture-induced negative affect. *PLoS*  
933 *biology*, 13(6):e1002180.
- 934 Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Machine learning*,  
935 20(3):273–297.
- 936 Decety, J. and Grèzes, J. (2006). The power of simulation: imagining one’s own and  
937 other’s behavior. *Brain research*, 1079(1):4–14.
- 938 Donahue, J., Anne Hendricks, L., Guadarrama, S., Rohrbach, M., Venugopalan, S.,  
939 Saenko, K., and Darrell, T. (2015). Long-term recurrent convolutional networks for  
940 visual recognition and description. In *Proceedings of the IEEE conference on*  
941 *computer vision and pattern recognition*, pages 2625–2634.
- 942 Fischl, B. (2012). Freesurfer. *Neuroimage*, 62(2):774–781.
- 943 Friston, K. J., Holmes, A. P., Worsley, K. J., Poline, J.-P., Frith, C. D., and Frackowiak,  
944 R. S. (1994). Statistical parametric maps in functional imaging: a general linear  
945 approach. *Human Brain Mapping*, 2(4):189–210.
- 946 Genovese, C. R., Lazar, N. A., & Nichols, T. (2002). Thresholding of statistical maps in  
947 functional neuroimaging using the false discovery rate. *Neuroimage*, 15(4), 870-878.
- 948 Glasser, M. F., Sotiropoulos, S. N., Wilson, J. A., Coalson, T. S., Fischl, B., Andersson,  
949 J. L., Xu, J., Jbabdi, S., Webster, M., Polimeni, J. R., et al. (2013). The minimal  
950 preprocessing pipelines for the human connectome project. *NeuroImage*, 80:105–  
951 124.
- 952 Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep  
953 feedforward neural networks. In *Proceedings of the 13th International Conference*  
954 *on Artificial Intelligence and Statistics (AISTATS)*, pages 249–256.
- 955 Glover, G. H. (1999). Deconvolution of impulse response in event-related bold fmri1.  
956 *NeuroImage*, 9(4):416–429.
- 957 Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. (2016). *Deep learning*,  
958 volume 1. MIT press Cambridge.
- 959 Gorgolewski, K., Burns, C. D., Madison, C., Clark, D., Halchenko, Y. O., Waskom, M.  
960 L., and Ghosh, S. S. (2011). Nipype: a flexible, lightweight and extensible  
961 neuroimaging data processing framework in python. *Frontiers in Neuroinformatics*,  
962 5:13.

- 963 Goutte, C. and Gaussier, E. (2005). A probabilistic interpretation of precision, recall  
964 and f-score, with implication for evaluation. In European Conference on Information  
965 Retrieval, pages 345–359. Springer.
- 966 Gramfort, A., Thirion, B., and Varoquaux, G. (2013). Identifying predictive regions  
967 from fmri with tv-l1 prior. In Pattern Recognition in Neuroimaging (PRNI), 2013  
968 International Workshop on, pages 17–20. IEEE.
- 969 Grosenick, L., Klingenberg, B., Katovich, K., Knutson, B., and Taylor, J. E. (2013).  
970 Interpretable whole-brain prediction analysis with graphnet. *NeuroImage*, 72:304–  
971 321.
- 972 Haxby, J. V., Gobbini, M. I., Furey, M. L., Ishai, A., Schouten, J. L., and Pietrini, P.  
973 (2001). Distributed and overlapping representations of faces and objects in ventral  
974 temporal cortex. *Science*, 293(5539):2425–2430.
- 975 Heekeren, H. R., Marrett, S., Bandettini, P. A., and Ungerleider, L. G. (2004). A general  
976 mechanism for perceptual decision-making in the human brain. *Nature*,  
977 431(7010):859.
- 978 Helfinstein, S. M., Schonberg, T., Congdon, E., Karlsgodt, K. H., Mumford, J. A., Sabb,  
979 F. W., Cannon, T. D., London, E. D., Bilder, R. M., and Poldrack, R. A. (2014).  
980 Predicting risky choices from brain activity patterns. *Proceedings of the National*  
981 *Academy of Sciences of the United States of America*, 111(7):2470–2475.
- 982 Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural*  
983 *Computation*, 9(8):1735–1780.
- 984 Holmes, A. and Friston, K. (1998). Generalisability, random effects & population  
985 inference. *NeuroImage*, 7:S754.
- 986 Hunt, L. T. and Hayden, B. Y. (2017). A distributed, hierarchical and recurrent  
987 framework for reward-based choice. *Nature Reviews Neuroscience*, 18(3):172.
- 988 Jang, H., Plis, S. M., Calhoun, V. D., and Lee, J.-H. (2017). Task-specific feature  
989 extraction and classification of fmri volumes using a deep neural network initialized  
990 with a deep belief network: Evaluation using sensorimotor tasks. *NeuroImage*,  
991 145:314–328.
- 992 Jorge, J., Van der Zwaag, W., and Figueiredo, P. (2014). Eeg–fmri integration for the  
993 study of human brain function. *Neuroimage*, 102:24–34.
- 994 Kiefer, J., Wolfowitz, J., et al. (1952). Stochastic estimation of the maximum of a  
995 regression function. *The Annals of Mathematical Statistics*, 23(3):462–466.
- 996 Kriegeskorte, N. and Douglas, P. K. (2018). Interpreting encoding and decoding  
997 models. arXiv preprint arXiv:1812.00278.

- 998 Kriegeskorte, N., Goebel, R., and Bandettini, P. (2006). Information-based functional  
999 brain mapping. *Proceedings of the National Academy of Sciences of the United*  
1000 *States of America*, 103(10):3863–3868.
- 1001 Lapuschkin, S., Binder, A., Montavon, G., Müller, K.-R., and Samek, W. (2016). The  
1002 layer-wise relevance propagation toolbox for artificial neural networks. *Journal of*  
1003 *Machine Learning Research*, 17(114):1–5.
- 1004 Lapuschkin, S., Wäldchen, S., Binder, A., Montavon, G., Samek, W., and Müller, K.-R.  
1005 (2019). Unmasking clever hans predictors and assessing what machines really learn.  
1006 *Nature Communications*, 10:1096.
- 1007 LeCun, Y., Bengio, Y., et al. (1995). Convolutional networks for images, speech, and  
1008 time series. *The handbook of brain theory and neural networks*, 3361(10):1995.
- 1009 LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *Nature*, 521(7553):436.
- 1010 Lindquist, M. A., Loh, J. M., Atlas, L. Y., and Wager, T. D. (2009). Modeling the  
1011 hemodynamic response function in fmri: efficiency, bias and mis-modeling.  
1012 *Neuroimage*, 45(1):S187– S198.
- 1013 Marban, A., Srinivasan, V., Samek, W., Fernández, J., and Casals, A. (2019). A  
1014 recurrent convolutional neural network approach for sensorless force estimation in  
1015 robotic surgery. *Biomedical Signal Processing and Control*, 50:134–150.
- 1016 McIntosh, A. R. and Lobaugh, N. J. (2004). Partial least squares analysis of  
1017 neuroimaging data: applications and advances. *Neuroimage*, 23:S250–S263.
- 1018 McLaughlin, N., Martinez del Rincon, J., and Miller, P. (2016). Recurrent convolutional  
1019 network for video-based person re-identification. In *Proceedings of the IEEE*  
1020 *conference on computer vision and pattern recognition*, pages 1325–1334.
- 1021 Mensch, A., Mairal, J., Thirion, B., and Varoquaux, G. (2018). Extracting universal  
1022 representations of cognition across brain-imaging studies. *arXiv preprint*  
1023 *arXiv:1809.06035*.
- 1024 Michel, V., Gramfort, A., Varoquaux, G., Eger, E., & Thirion, B. (2011). Total variation  
1025 regularization for fMRI-based prediction of behavior. *IEEE transactions on medical*  
1026 *imaging*, 30(7), 1328-1340.
- 1027 Montavon, G., Bach, S., Binder, A., Samek, W., and Müller, K.-R. (2017). Explaining  
1028 nonlinear classification decisions with deep taylor decomposition. *Pattern*  
1029 *Recognition*, 65:211–222.
- 1030 Montavon, G., Samek, W., and Müller, K.-R. (2018). Methods for interpreting and  
1031 understanding deep neural networks. *Digital Signal Processing*, 73:1–15.

- 1032 Müller, K.-R., Mika, S., Rätsch, G., Tsuda, K., and Schölkopf, B. (2001). An  
1033 introduction to kernel-based learning algorithms. *IEEE transactions on neural*  
1034 *networks*, 12(2):181–201.
- 1035 Naselaris, T., Kay, K. N., Nishimoto, S., and Gallant, J. L. (2011). Encoding and  
1036 decoding in fmri. *Neuroimage*, 56(2):400–410.
- 1037 Nie, D., Zhang, H., Adeli, E., Liu, L., and Shen, D. (2016). 3d deep learning for multi-  
1038 modal imaging-guided survival time prediction of brain tumor patients. In  
1039 *International Conference on Medical Image Computing and Computer-Assisted*  
1040 *Intervention*, pages 212–220. Springer.
- 1041 Olson, I. R., Plotzker, A., and Ezzyat, Y. (2007). The enigmatic temporal pole: a review  
1042 of findings on social and emotional processing. *Brain*, 130(7):1718–1731.
- 1043 Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent  
1044 neural networks. In *International Conference on Machine Learning*, pages 1310–  
1045 1318.
- 1046 Petrov, D., Kuznetsov, B. A., van Erp, T. G., Turner, J. A., Schmaal, L., Veltman, D.,  
1047 Wang, L., Alpert, K., Isaev, D., Zavaliangos-Petropulu, A., et al. (2018). Deep  
1048 learning for quality control of subcortical brain 3d shape models. *arXiv preprint*  
1049 *arXiv:1808.10315*.
- 1050 Plis, S. M., Hjelm, D. R., Salakhutdinov, R., Allen, E. A., Bockholt, H. J., Long, J. D.,  
1051 Johnson, H. J., Paulsen, J. S., Turner, J. A., and Calhoun, V. D. (2014). Deep  
1052 learning for neuroimaging: a validation study. *Frontiers in neuroscience*, 8:229.
- 1053 Poldrack, R. A., Barch, D. M., Mitchell, J., Wager, T., Wagner, A. D., Devlin, J. T.,  
1054 Cumba, C., Koyejo, O., and Milham, M. (2013). Toward open sharing of task-based  
1055 fmri data: the openfmri project. *Frontiers in Neuroinformatics*, 7:12.
- 1056 Reverberi, C., Kuhlen, A. K., Seyed-Allaei, S., Greulich, R. S., Costa, A., Abu- talebi,  
1057 J., and Haynes, J.-D. (2018). The neural basis of free language choice in bilingual  
1058 speakers: Disentangling language choice and language execution. *NeuroImage*,  
1059 177:108–116.
- 1060 Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning represen-tations  
1061 by back-propagating errors. *nature*, 323(6088):533.
- 1062 Ryali, S., Supekar, K., Abrams, D. A., and Menon, V. (2010). Sparse logistic regression  
1063 for whole-brain classification of fmri data. *NeuroImage*, 51(2):752–764.
- 1064 Samek, W., Wiegand, T., and Müller, K.-R. (2018). Explainable artificial intelligence:  
1065 Understanding, visualizing and interpreting deep learning models. *ITU Journal: ICT*  
1066 *Discoveries - Special Issue 1 - The Impact of Artificial Intelligence (AI) on*  
1067 *Communication Networks and Services*, 1(1):39–48.

- 1068 Sarraf, S. and Tofighi, G. (2016). Classification of alzheimer's disease using fmri data  
1069 and deep learning convolutional neural networks. arXiv preprint arXiv:1603.08631.
- 1070 Schölkopf, B. and Smola, A. J. (2002). Learning with kernels: support vector machines,  
1071 regularization, optimization, and beyond. MIT press.
- 1072 Schuck, N. W., Cai, M. B., Wilson, R. C., and Niv, Y. (2016). Human orbitofrontal  
1073 cortex represents a cognitive map of state space. *Neuron*, 91(6):1402–1412.
- 1074 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014).  
1075 Dropout: A simple way to prevent neural networks from overfitting. *The Journal of*  
1076 *Machine Learning Research*, 15(1):1929–1958.
- 1077 Sturm, I., Lapuschkin, S., Samek, W., and Müller, K.-R. (2016). Interpretable deep  
1078 neural networks for single-trial eeg classification. *Journal of neuroscience methods*,  
1079 274:141–145.
- 1080 Suk, H.-I., Lee, S.-W., Shen, D., Initiative, A. D. N., et al. (2014). Hierarchical feature  
1081 representation and multimodal fusion with deep learning for ad/mci diagnosis.  
1082 *NeuroImage*, 101:569–582.
- 1083 Thomas, A. W., Müller, K. R., & Samek, W. (2019). Deep Transfer Learning for  
1084 Whole-Brain FMRI Analyses. In *OR 2.0 Context-Aware Operating Theaters and*  
1085 *Machine Learning in Clinical Neuroimaging* (pp. 59-67). Springer, Cham.
- 1086 Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the*  
1087 *Royal Statistical Society. Series B (Methodological)*, pages 267–288.
- 1088 Tikhonov, A. N. (1943). On the stability of inverse problems. In *Dokl. Akad. Nauk*  
1089 *SSSR*, volume 39, pages 195–198.
- 1090 Uğurbil, K., Xu, J., Auerbach, E. J., Moeller, S., Vu, A. T., Duarte-Carvajalino, J. M.,  
1091 Lenglet, C., Wu, X., Schmitter, S., Van de Moortele, P. F., et al. (2013). Pushing  
1092 spatial and temporal resolution for functional and diffusion mri in the human  
1093 connectome project. *NeuroImage*, 80:80–104.
- 1094 Van Essen, D. C., Smith, S. M., Barch, D. M., Behrens, T. E., Yacoub, E., Ugurbil, K.,  
1095 Consortium, W.-M. H., et al. (2013). The wu-minn human connectome project: an  
1096 overview. *NeuroImage*, 80:62–79.
- 1097 Wager, T. D., Atlas, L. Y., Lindquist, M. A., Roy, M., Woo, C.-W., and Kross, E.  
1098 (2013). An fmri-based neurologic signature of physical pain. *New England Journal*  
1099 *of Medicine*, 368(15):1388–1397.
- 1100 Weygandt, M., Schaefer, A., Schienle, A., and Haynes, J.-D. (2012). Diagnosing  
1101 different binge-eating disorders based on reward-related brain activation patterns.  
1102 *Human brain mapping*, 33(9):2135–2146.

1103 Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C., and Wager, T. D. (2011).  
1104 Large-scale automated synthesis of human functional neuroimaging data. *Nature*  
1105 *Methods*, 8(8):665.

1106 Yousefnezhad, M. and Zhang, D. (2018). Anatomical pattern analysis for decoding  
1107 visual stimuli in human brains. *Cognitive Computation*, 10(2):284–295.

1108

## 1109 **Figure legends**

1110 Figure 1: Illustration of the DeepLight approach. A whole-brain fMRI volume is sliced  
1111 into a sequence of axial images. These images are then passed to a DL model consisting  
1112 of a convolutional feature extractor, an LSTM and an output unit. First, the  
1113 convolutional feature extractor reduces the dimensionality of the axial brain slices  
1114 through a sequence of eight convolution layers. The resulting sequence of higher-level  
1115 slice representations is then fed to a bi-directional LSTM, modeling the spatial  
1116 dependencies of brain activity within and across brain slices. Lastly, the DL model  
1117 outputs a decoding decision about the cognitive state underlying the fMRI volume,  
1118 through a softmax output layer with one output neuron per cognitive state in the data.  
1119 Once the prediction is made, DeepLight utilizes the LRP method to decompose the  
1120 prediction into the contributions (or relevance) of the single input voxels to the  
1121 prediction. Thereby, enabling an analysis of the association between fMRI data and  
1122 cognitive state.

1123 Figure 2: Group-level decoding performance of DeepLight, the searchlight analysis and  
1124 whole-brain lasso. A: Confusion matrix of DeepLight's decoding decisions. B: Average  
1125 decoding performance of DeepLight over the course of an experiment block. C-D:  
1126 Confusion matrix for the decoding decisions of the group-level searchlight analysis (C)  
1127 and whole-brain lasso (D). E: Average decoding accuracy of the searchlight (green),  
1128 whole-brain lasso (blue) and DeepLight (red), when these are repeatedly trained on a  
1129 subset of the subjects from the full training dataset. Black dashed horizontal lines  
1130 indicate chance level.

1131 Figure 3: Subject-level decoding performance comparison of DeepLight (red) to the  
1132 searchlight analysis (A; green) and whole-brain lasso (B; blue). Black scatter points  
1133 indicate the average decoding accuracy for a subject. Colored lines indicate the average  
1134 decoding accuracy across all 30 test subjects.

1135 Figure 4: Group-level brain maps for each cognitive state and analysis approach: A:  
1136 Results of a NeuroSynth meta-analysis for the terms "body", "face", "place" and "tools".  
1137 The brain maps were thresholded at an expected false discovery rate of 0.01. Red boxes  
1138 highlight the regions-of-interest for each cognitive state. B: Results of the GLM group-  
1139 level analysis. The brain maps of the GLM analysis were thresholded at an expected  
1140 false discovery rate of 0.1. C-E: Results of the group-level searchlight analysis (C),  
1141 whole-brain lasso (D) and DeepLight (E). The brain maps of the searchlight analysis,  
1142 whole-brain lasso, and DeepLight were thresholded at the 90th percentile of their  
1143 values. Note that the values of the brain maps are on different scales between analysis

1144 approaches, due to their different statistical nature. All brain maps are projected onto the  
1145 inflated cortical surface of the FsAverage5 surface template (Fischl, 2012).

1146 Figure 5: Exemplary DeepLight brain maps for each of the four cognitive states on  
1147 different levels of data granularity for a single subject. All brain maps belong to the  
1148 subject with the highest decoding accuracy in the held-out test dataset. A: Average  
1149 relevance maps for all correctly classified TRs of the subject (with an average of 47 TRs  
1150 per cognitive state). B: Average relevance maps for all correctly classified TRs of the  
1151 first experiment block of each cognitive state in the first experiment run (with an  
1152 average of 12 TRs per cognitive state). C: Exemplar relevance maps for a single TR of  
1153 the first experiment block of each cognitive state in the first experiment run. All  
1154 relevance maps were thresholded at the 90th percentile of their values and projected  
1155 onto the inflated cortical surface of the FsAverage5 surface template (Fischl, 2012).

1156 Figure 6: DeepLight analysis of the temporo-spatial distribution of brain activity in the  
1157 first experiment block of the face and place stimulus classes in the second experiment  
1158 run of the held-out test dataset. A-B: Average predicted probability that the fMRI data  
1159 collected at each sampling time point belongs to each of the four cognitive states. C &  
1160 E: Results of a meta-analysis with the NeuroSynth database for the face and place  
1161 stimulus classes (for details on the meta-analysis, see Supplementary Information  
1162 Section 1). D & F: Group-level brain maps for seven fMRI sampling time points from  
1163 the experiment block. Each group-level brain map at each time point is computed as an  
1164 average over the relevance maps of each subject for this time point. Each group-level  
1165 brain map is thresholded at the 90th percentile of its values. All brain maps are  
1166 projected onto the inflated cortical surface of the FsAverage5 surface template (Fischl,  
1167 2012). G-H: F1-score for each group-level brain map at each sampling time point of the  
1168 experiment block. The F1-score quantifies the similarity between the group-level brain  
1169 map and the results of the meta-analysis (C & E) (for further details on the F1-score, see  
1170 Section 3.2.2 and Supplementary Information Section 2). Red indicates the results of the  
1171 F1-score comparison for the brain maps of DeepLight, whereas blue indicates the results  
1172 of this comparison for the brain maps of the whole-brain lasso analysis (for further  
1173 details on the F1-comparison for the whole-brain lasso analysis, see Section 3.4).

Figure 1.JPEG

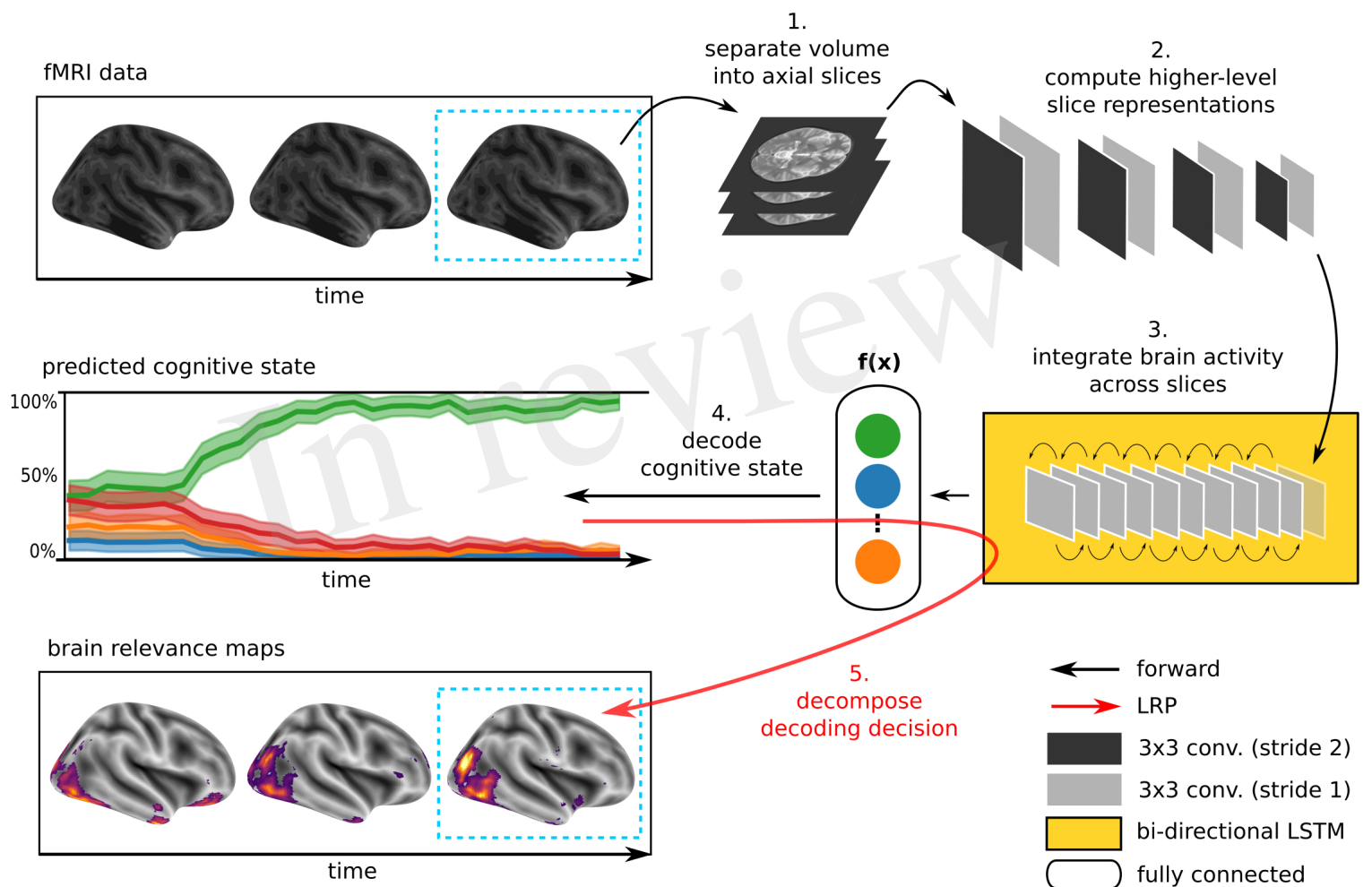




Figure 2.JPEG

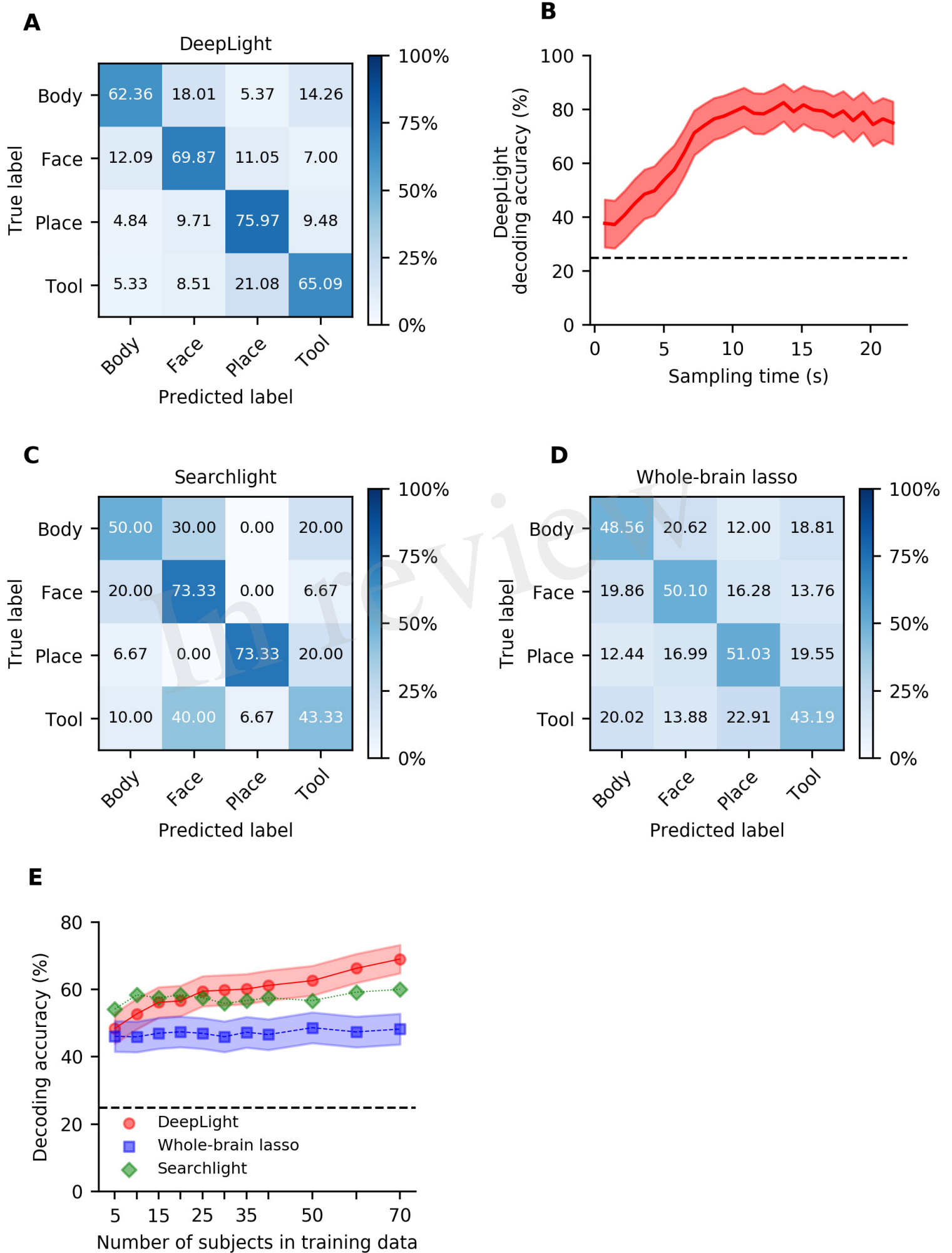


Figure 3.JPEG

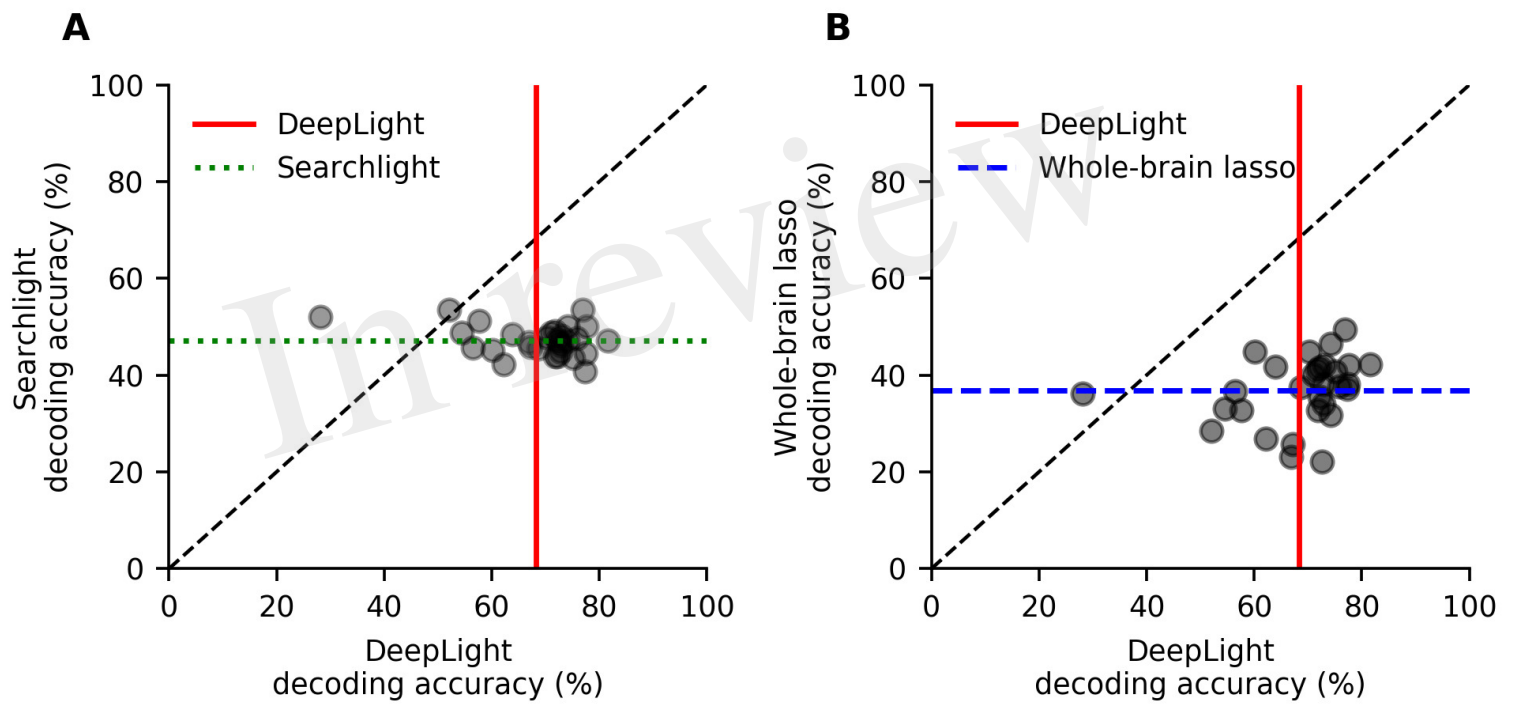


Figure 4.JPEG

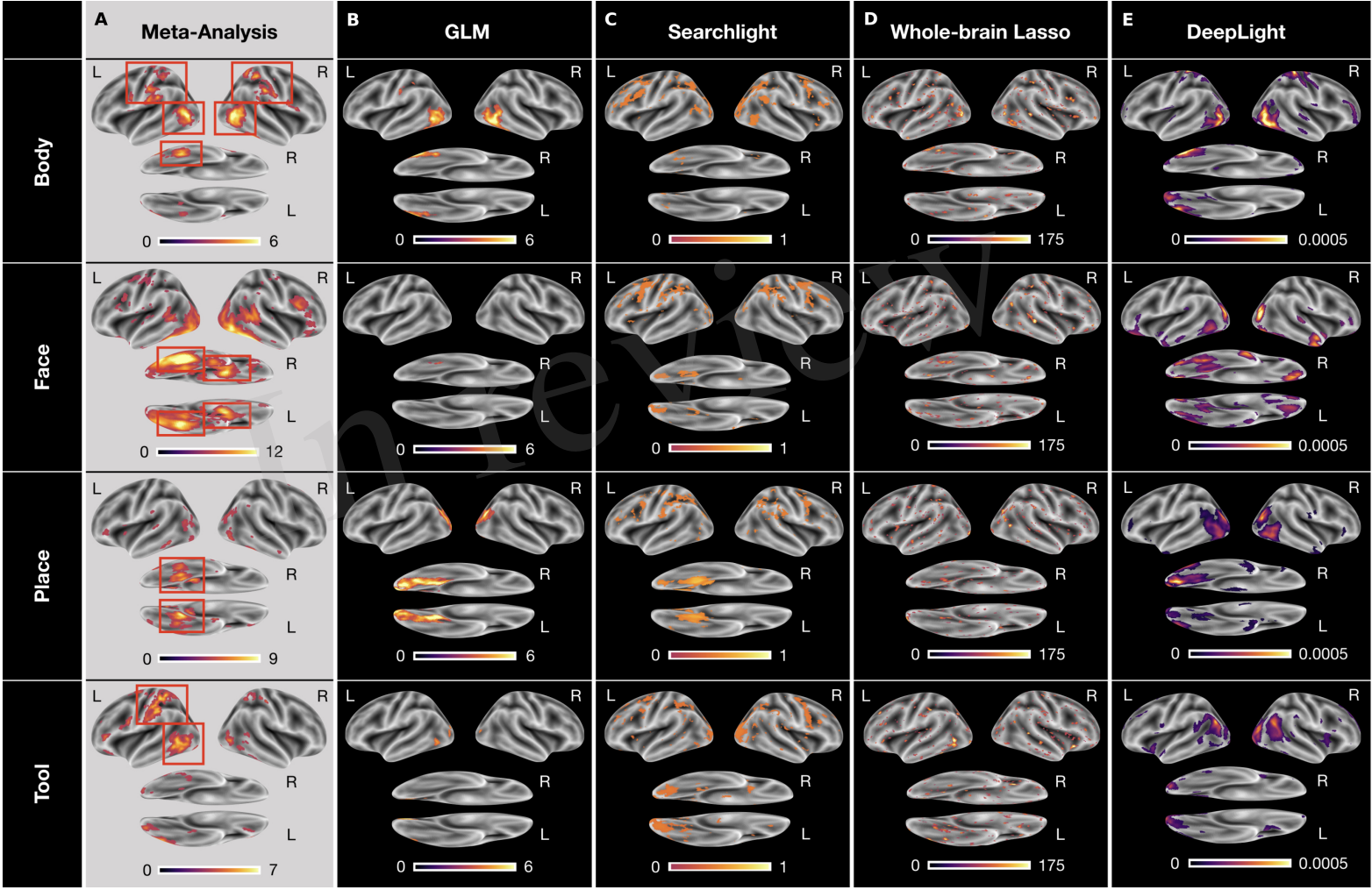


Figure 5.JPEG

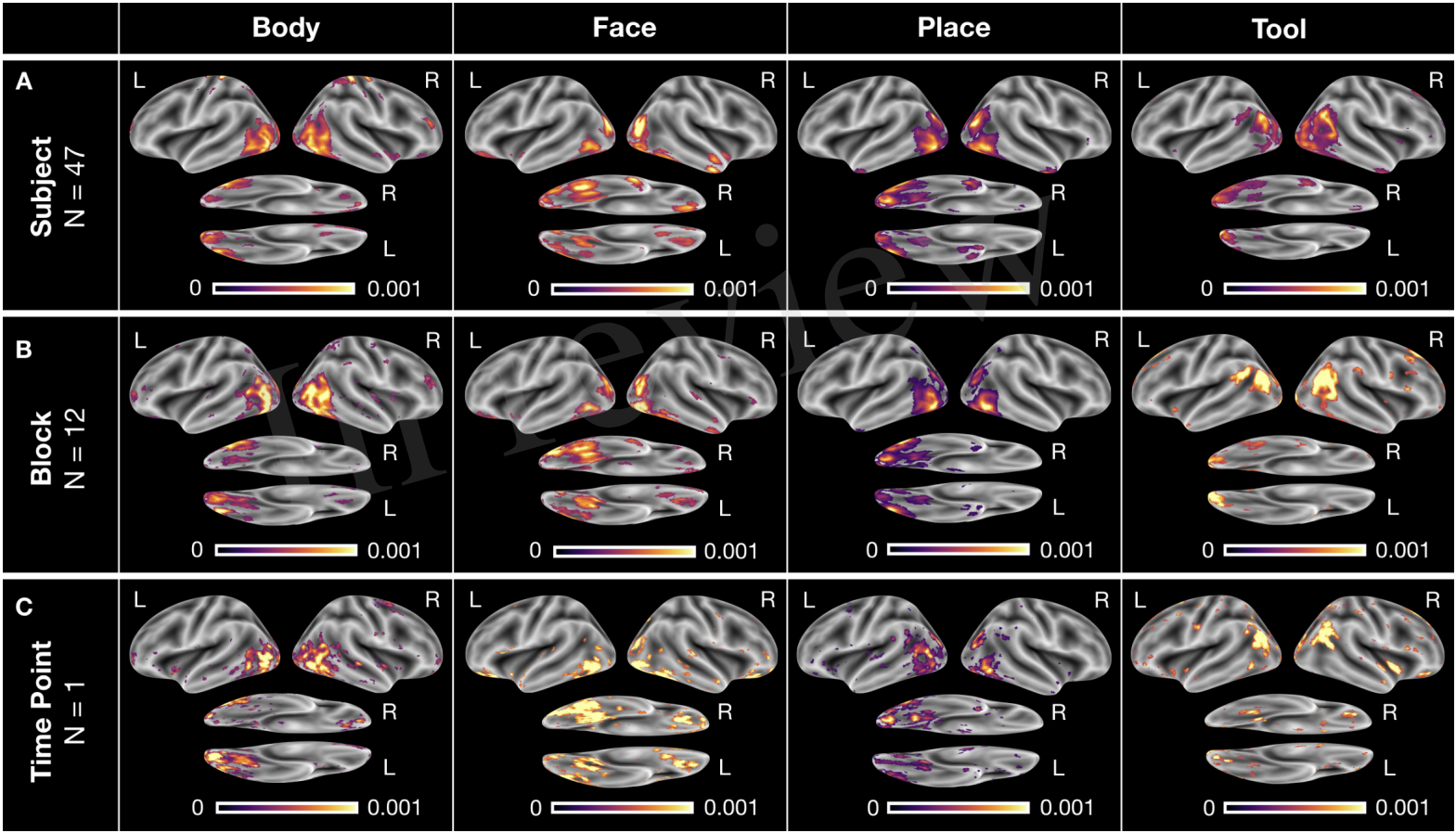


Figure 6.JPEG

