REPRESENTATIONS IN A DEEP END-TO-END DRIVING MODEL PREDICT HUMAN BRAIN ACTIVITY IN AN ACTIVE DRIVING TASK

Anonymous authors

Paper under double-blind review

ABSTRACT

Understanding how cognition and learned representations give rise to intelligent behavior is a fundamental goal in both machine learning and neuroscience. However, in both domains, the most well-understood behaviors are passive and openloop, such as image recognition or speech processing. In this work, we compare human brain activity measured via functional magnetic resonance imaging with deep neural network (DNN) activations for an active taxi-driving task in a naturalistic simulated environment. To do so, we used DNN activations to build voxelwise encoding models for brain activity. Results show that encoding models for DNN activations explain significant amounts of variance in brain activity across many regions of the brain. Furthermore, each functional module in the DNN explains brain activity in a distinct network of functional regions in the brain. The functions of each DNN module correspond well to the known functional properties of its corresponding brain regions, suggesting that both the DNN and the human brain may partition the task in a similar manner. These results represent a first step towards understanding how humans and current deep learning methods agree or differ in active closed-loop tasks such as driving.

025

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

031 Fully autonomous vehicles now share the roads with human drivers in several major cities. In 032 contrast with most previous commercially successful applications of robotics, such as robotized 033 warehouses (Azadeh et al., 2019), drone medication delivery in rural areas (Demuyakor, 2020), and 034 infrastructure inspection (Lattanzi & Miller, 2017), autonomous driving systems perform the same tasks as humans do in diverse and complex environments involving dynamic interactions with other 035 agents. Ensuring that the algorithms used in autonomous driving are effective and safe to very low margins of error poses a significant challenge. Furthermore, most autonomous driving systems 037 use deep neural networks (DNNs). DNNs are notoriously difficult to analyze and interpret (Zhang et al., 2021), which makes verification of these desirable properties an additional problem. However, there are also exciting opportunities: for the first time, human and artificial intelligences (AIs) are 040 performing the same complex sensorimotor and social tasks in a shared environment. This raises 041 many interesting questions at the intersection of artificial intelligence and human neuroscience: to 042 what extent do the human brain and DNNs rely on similar representations when integrating diverse 043 sensory information to produce task-appropriate behavior? Do the human brain and DNN share 044 any functional organization for solving complex, closed-loop tasks? Do they make similar kinds of 045 inferences about each other when engaged in a shared task environment?

In this work, we provide some initial insights into these questions by comparing two sets of features:
brain activity from human drivers recorded with functional magnetic resonance imaging (fMRI), and
activations from the Learning from All Vehicles (LAV) driving DNN (Chen & Krähenbühl, 2022)
when given similar inputs to those seen by the human drivers. To perform this comparison, we
regress the LAV activations onto the brain activity to build a voxelwise encoding model (Naselaris
et al., 2011; Nunez-Elizalde et al., 2019), then test this model's predictive power on held-out brain
data. We find that the LAV voxelwise encoding model explains a substantial amount of variance in
human brain activity. Furthermore, features from the modules in the LAV DNN explain variance in
distinct functional regions in the brain. The known functional properties of these regions align well



071 Figure 1: Using deep neural network activations to model brain activity during a naturalistic simulated driving task. Human subjects drove through a large virtual city in the CARLA simulator and performed a taxi-driver task while brain activity was recorded with functional magnetic reso-073 nance imaging (fMRI). The recorded experiment state was used to generate the image and LiDAR 074 sensor inputs that were then passed to Learning from All Vehicles (LAV), a deep neural network 075 driving model. Activations in the LAV model in response to these inputs were regressed onto brain 076 activity to produce a predictive encoding model of brain activity during the driving task. This model 077 can be used to investigate and quantify the alignment of the neural network driving model to the 078 representations the human brain employs during driving. 079

with the corresponding modules. Our results are an exciting first step towards investigating the cognitive and representational basis for human and AI driving by leveraging each to better understand the other.

2 RELATED WORK

2.1 HUMAN BRAIN ACTIVITY DURING DRIVING

The cognitive processes underlying human driving behavior have been studied in many experiments 090 and with a variety of brain imaging modalities, including fMRI, electroencephalography (EEG), 091 functional near-infrared spectroscopy (fNIRS), and magnetoencephalography (MEG) (Haghani et al., 2021). Of these modalities, fMRI provides the highest spatial resolution and enables the 092 localization of brain activity (Gross, 2019), and has been used to localize various aspects of driving to functional brain regions. Both Schweizer et al. (2013) and Spiers & Maguire (2007) have used 094 fMRI to examine brain activity from subjects performing simulated driving tasks. Navarro et al. 095 (2018) pooled these and seven other fMRI studies of driving to find characteristic patterns of brain 096 activity for low-level tasks such as turning, medium-level tasks such as planning to overtake, and high-level tasks such as route planning, and related these patterns to conceptual (not quantitative or 098 predictive) models of human driving. Unfortunately, these studies only used contrast-based fMRI 099 analysis methods to localize driving-related activity to certain brain regions, and did not investigate 100 hypotheses about specific representations and computations in the brain during driving. Also, these 101 studies provide no insights into AI driving algorithms.

102

081

082

083

084 085

087

088

103 2.2 DNNs as encoding models of brain activity

Encoding models are a powerful method for the analysis of brain activity. In this approach, the time series of brain activity is modelled as a function of stimulus and task feature time series. Encoding models have been extensively used to study how the brain represents visual inputs. Features used to fit encoding models of brain activity in these studies range from simple transformations of the visual

input, such as Gabor wavelets (Carandini et al., 2005), to deep neural networks trained on vision tasks (Agrawal et al., 2014; Güçlü & Van Gerven, 2015; Takagi & Nishimoto, 2023). Because previous studies using DNN-based encoding models have focused on visual processing, they used data from tasks in which subjects observed stimuli passively. Here we aim to apply the encoding model approach to data from an active sensorimotor task in a naturalistic environment. While vision is an important part of this task, active driving also involves planning, motor control, and social interaction with other agents.

- 115
- 116 117

3 VOXELWISE MODELING WITH DEEP NEURAL NETWORKS FOR DRIVING

To compare deep neural network activations with brain activity, we first collected a dataset of driving behavior and brain activity from human subjects. We then used the recorded behavior from this dataset to generate a time series of appropriate inputs to the LAV driving network. Next, we gave these inputs to the LAV driving network, and extracted features representing the activity of the LAV network. Finally, we built voxelwise models of brain activity by regressing the LAV driving network features against the brain activity.

124

126

125

3.1 MEASURING BRAIN ACTIVITY WITH FMRI

We used functional magnetic resonance imaging (fMRI) to record brain activity from three sub-127 jects performing a taxi-driver task in a simulator. We used Unreal Engine 4 and the CARLA plugin 128 to implement a driving simulator that contains a large 2×3 km virtual city populated by dynamic 129 pedestrians and vehicles. Prior to recording, subjects learned the layout of the world. During record-130 ing, subjects controlled a virtual car with an MR-compatible steering wheel and pedals (illustrated 131 in Fig.1), performing the taxi-driver task, in which they were cued to navigate to particular loca-132 tions. Blood-oxygenation-level dependent (BOLD) (Ogawa et al., 1990) activity from the brain was 133 recorded by the MRI scanner as subjects performed the task. Data were collected in 11-minute runs (180 minutes for subjects 1 and 2, and 110 minutes for subject 3) at a temporal resolution of 2.0045 134 seconds and a voxel (3D pixel) size of $2.24 \times 2.24 \times 3.5$ mm³ (matrix size = 100×100 voxels, 30 135 axial slices). Please refer to appendix A.3 for more details on data collection and preprocessing. The 136 experimental procedures were approved by the Institutional Review Board at [redacted], and written 137 informed consent was obtained from all subjects. 138

139 140

3.2 Measuring deep neural network activity

141 3.2.1 DNN DRIVING MODEL

143 For this work, we chose to examine a pre-trained driving deep neural network developed by Chen & Krähenbühl (2022), Learning from All Vehicles (LAV). LAV was developed for the sensors track of 144 the CARLA Leaderboard challenge (version 1.0), where teams compete to submit algorithms that 145 can complete a sequence of driving scenarios while obeying traffic rules and avoiding collisions. 146 LAV achieved first place in the 2021 edition of the challenge, with a route completion rate of 94% 147 (Chen & Krähenbühl, 2022). It uses an imitation learning-based training objective to mimic "expert" 148 driving behavior provided by the built-in CARLA AI driver. The CARLA AI driver uses ground-149 truth simulator state to generate appropriate driving actions, and LAV attempts to imitate its behavior 150 with only RGB images and LiDAR as inputs. LAV was trained on CARLA version 0.9.10, which is 151 slightly different than the modified version of CARLA we used in our human subject experiments, 152 but we verified that LAV performed reasonable inferences when run with inputs from our simulator 153 (see Fig.2a for example outputs on our data).

The LAV DNN consists of an end-to-end architecture, but is split into different distinct modules that each perform a specific sub-task for driving (encouraged via extra module-specific objective functions during training). We used this architecture to conceptually group the measured DNN activations by the sub-task with which they are associated. In addition to processing the RGB and LiDAR inputs and producing control outputs, these modules also include "intermediate" representations such as a bird's-eye-view of the local environment. The modules are as follows:

- 160 161
- 1. **Semantic segmentation:** This module predicts the semantic class of visual features in RGB images from the vehicle's cameras. It receives three images as input, one from a

camera facing forwards and two at yaw offsets of \pm 60 degrees. It categorizes each pixel in these images as one of 5 semantic classes: roads, road markings, vehicles, pedestrians, and other. The architecture is an ERFNet (Romera et al., 2017), a type of convolutional neural network specialized for segmentation.

- 2. Bird's-eye-view (BEV) perception: This module predicts local scene features, including 167 the layout of the road and the locations of other vehicles, in a BEV reference frame fixed on the ego vehicle. It receives 3D LiDAR points as input, where points that lie within the field 169 of view of the RGB cameras are "painted" with their semantic class label from the semantic 170 segmentation module. The architecture is split into two learned components: a point pillar 171 net and a perception backbone. The point pillar splits painted LiDAR points into a 320 x 320 grid of the environment surrounding the ego vehicle as seen from above, and returns 172 a feature vector for each bin of the grid. (Since the point pillar network operates on an 173 unordered and variable-length set of LiDAR points, we only include the final grid-based feature vector in the features extracted from the point pillar.) The perception backbone 175 then passes these grid features through a series of convolutional layers to produce a feature 176 map, which is decoded into a semantic segmentation of the local road layout as well as the 177 position and orientation of other vehicles, all of which are in a top-down, bird's-eye-view perspective.
- 179 3. **Planning:** This module generates a trajectory plan for the ego vehicle. It takes as input the BEV feature map from the previous module, cropped and centered on the ego vehicle. 181 This feature map is first passed through a ResNet-18 (He et al., 2016), which performs 182 driving plan-agnostic processing. Then, for each of 6 possible high-level commands (turn 183 left / right, go straight through an intersection, stay in the current lane, and lane change left / right), a gated recurrent unit (GRU) network predicts the future position in 0.25 s intervals (up to a planning horizon of 2.5 s). Finally, a second GRU network is used to 186 refine the intermediate set of trajectories into a final plan conditioned on a target waypoint indicating the desired route. Since the human drivers in our experiment are not following a 187 fixed route, we generate target waypoints for the planning module using the position of the 188 human driver 40 m further along their future trajectory. 189
- 4. **Trajectory prediction:** This module generates predicted trajectories for each nearby detected vehicle. The inputs and architecture are exactly the same as in the planning module, 192 except that the BEV feature map is cropped and centered around the detected vehicle, and there is no target waypoint (which is unknown for the other vehicles) and therefore no final trajectory-refinement GRU. Since the true high-level command is also unknown for the de-194 tected vehicles, this module also uses a single linear layer to predict the probability of each possible command, resulting in a probability distribution over the trajectories associated 196 with each command.
- 5. Hazard detection: This module predicts whether there is a hazard present that necessitates braking. It receives four RGB images as input (the same inputs as the semantic segmentation module, plus an additional front-facing image at a higher zoom level). The architecture 200 is a ResNet-18 followed by a linear layer for classification. During training, the output of the ResNet-18 network is additionally passed into another network which outputs a seman-202 tic segmentation, but this is not part of inference and not included in our extracted features.
- 6. **Control:** This module takes the outputs of the planning, prediction, and hazard detection 204 modules, and determines the final control action that should be applied. Specifically, if the 205 hazard detection model predicts that braking needs to occur or if the planned trajectory for 206 the ego vehicle intersects the predicted trajectory of another vehicle, the control module will 207 brake to stop the car. If the control module decides that no braking is needed, it uses PID 208 control to calculate an acceleration and steering command to follow the planned trajectory provided by the planning module. Unlike the other modules, the controller is not a DNN 210 and has no learned parameters.
- 211

201

203

162

163

164

165

212 3.2.2 MEASURING DNN ACTIVITY 213

To measure the activity of the network, we used recorded simulator state data from the human 214 subjects' driving sessions to generate RGB images and 3D LiDAR scans at 15 frames per second 215 (FPS). There are four cameras, two pointing forwards (one with a smaller field of view and longer focal length) and two pointing at an offset angle of 60 degrees to the left and right of the forward
cameras. Because LAV does not have any form of memory other than selecting the current target
waypoint from a pre-defined route, we used the recorded human data as the input at every frame, and
ignored LAV's final control outputs without disrupting the operation of the LAV network. Doing so
enabled us to measure the activation of the LAV network in response to the same simulator state
encountered by the subjects and align the activity of the LAV DNN and the human brain activity for
our models.

223 Across all modules, the LAV DNN contains on the order of 10^8 measurable values that vary with the 224 input data. This large number makes it computationally intractable to store the activations from all 225 units and regress them against brain activity. To overcome this computational challenge and reduce 226 the dimensionality of the LAV network activations, for each module except the controller, we used a sparse random projection (SRP) of k components per module (Li et al., 2006), using the imple-227 mentation in Scikit-learn to generate the random projection matrix (Pedregosa et al., 2011). SRPs 228 approximately preserve the distance between points with high probability and dramatically speed up 229 linear regression (Woodruff et al., 2014). We provide a proof in appendix A.1 that solving ridge re-230 gression with features projected with an SRP results in an approximate solution to the original ridge 231 regression problem with high probability, where choosing k large enough results in low approxima-232 tion error and high success probability. We used the same k for each module because it controls for 233 the number of regression features when comparing encoding model performances across modules. 234 In this work, we found k = 20,000 to produce stable regression results across multiple random 235 projections. 236

For the controller module, we used the planned trajectory and brake probability for a total of 121 features, and therefore did not apply any dimensionality reduction. The planned trajectory is in a BEV reference frame.

- 240
- 241
- 242 243

3.3 VOXELWISE MODELING

244 245 246

We applied the voxelwise modeling (VM; Naselaris et al. (2011); Huth et al. (2012; 2016)) frame-247 work to the features extracted from LAV to determine whether a DNN and the human brain rely 248 on similar representations while driving, and whether the DNN and the human brain share any 249 functional organization for active driving. In VM, the time series of activity in each brain voxel is 250 modeled as a linear combination of the time series of all features. Models are evaluated by predict-251 ing brain activity on a held-out dataset. High prediction performance by a set of features, or feature 252 space, suggests that the brain represents information as parameterized by that feature space (Nase-253 laris et al., 2011). VM has enabled many new insights about brain activity, such as in experiments 254 where subjects listen to narrative stories (Huth et al., 2016; Deniz et al., 2019) or watch movies (Nishimoto et al., 2011; Huth et al., 2012). 255

256 In this work, we treat each module of the LAV network as a separate feature space. To capture the 257 contributions of each module to the final predictive performance of the encoding model, we used 258 banded ridge regression (Nunez-Elizalde et al., 2019; Dupré la Tour et al., 2022), which imposes a 259 different ridge regularization parameter on each feature space. After regression, the overall model performance was quantified by computing the R^2 (explained variance) between predicted and actual 260 activity in each voxel on a held-out test set, and individual model performances for each of the fea-261 ture spaces were determined by partitioning the R^2 between them (Pratt, 1987; Dupré la Tour et al., 262 2022). Intuitively, each of the feature spaces can be considered as a hypothesis of how information 263 may be represented in (some part of) the brain; the split R^2 for each feature space provides a test for 264 its corresponding hypothesis. 265

The BOLD signal captured by fMRI does not respond to stimuli instantaneously but rather over a period of several seconds (the hemodynamic response function). We modeled this by implementing
a finite impulse response (FIR) filter with four delays of 2 s for each LAV feature, capturing up to
s in the past. The best-performing feature weights, regularization parameters, and FIR filter shape were empirically selected by cross-validating over 10,000 random samples.



Figure 2: Features from different modules in a deep neural network driving model predict 307 activity in distinct functional brain networks. a): the architecture of the Learning from All Vehi-308 cles (LAV) neural network. Each colored oval represents a different module, with their inputs and 309 outputs shown with rectangles. Activations from each module in response to the stimuli seen by 310 human subjects were used to fit voxelwise encoding models to the recorded brain activity. Model 311 performance was quantified as the explained variance (R^2) in brain activity on a held-out dataset 312 not used during model fitting. b): A group-level map of the best performing LAV modules across 313 the flattened cortical surface. Voxels with a p-value of < 0.01 in at least 2 subjects are shown in the 314 color corresponding to the module with the highest explained variance. Features from the seman-315 tic segmentation and hazard detection networks, which process RGB images, best explain variance 316 in visual areas V1-V4. The planning module best explains variance in the sensorimotor areas and 317 the intraparietal sulcus (IPS). The control module has a similar pattern of explained variance to the planning module but additionally provides a good fit to the retrosplenial cortex (RSC) and per-318 ahippocampal place area (PPA). The bird's-eye-view (BEV) perception module and the planning 319 module for other vehicles do not explain variance in consistently localized regions across subjects. 320 c): Group-level encoding model performance for the semantic segmentation, hazard detection, plan-321 ning, and control modules across all subjects. The consistent patterns of brain activity explained by 322 each of these modules suggests that humans and the driving network share an analogous organiza-323 tion of cognitive information in the active driving task.

³²⁴ 4 RESULTS AND DISCUSSION

325 326

Our encoding model explains variance across much of the brain, including much of the visual and motor cortices, as illustrated in Fig.2. The model also explains brain activity in parts of the prefrontal cortex. In contrast, passive vision-only experiments only engage visual perceptual areas, and encoding models for these experiments are limited to explaining variance in these areas (Yamins et al., 2014; la Tour et al., 2021; Takagi & Nishimoto, 2023).

331 To determine whether the driving DNN shares any functional organizational principles with the 332 brain, we used the split R^2 scores to identify the best-performing driving DNN module for every 333 voxel (Fig.2b). Quantitatively, we find that the spatial organization of the best-performing modules 334 for each voxel is non-random in all three subjects (see appendix A.3 for more details on statistical 335 tests). Qualitatively, We find that earlier LAV modules, such as the semantic segmentation module, 336 best explain brain activity in visual areas, whereas later modules, such as the planning module, best explain brain activity in sensorimotor and higher-level cognitive brain regions. Overall, this result 337 suggests that the representations learned by the driving DNN may be similar to those used by the 338 human brain, despite significant differences in perceptual inputs (human drivers have no LiDAR-like 339 active sensing perception modalities, and our subjects were able to see only the image from a single 340 camera with a smaller field of view than the four cameras used by LAV), training data (humans learn 341 to drive primarily through closed-loop control rather than imitation learning), and physical substrate 342 (real neurons behave very differently from artificial neural network units) between the two systems. 343 In the following sections, we describe the predictive performance of each driving network module 344 in more detail.

345 346

347

4.1 SEMANTIC SEGMENTATION AND HAZARD DETECTION

348 The semantic segmentation and hazard detection modules, which both process RGB images from 349 the vehicle cameras, explain brain activity almost exclusively in the posterior visual cortex. This pattern of explained variance is qualitatively similar to previous encoding model studies that used 350 features derived from convolutional neural networks (CNNs) trained on vision tasks such as image 351 recognition (la Tour et al., 2021; Güçlü & Van Gerven, 2015). The hazard detection network explains 352 more variance in the inferior part of the visual cortex, which corresponds to the upper half of the 353 visual field. This bias may be due to the fact that, when approaching a vehicle, participants are 354 likely to shift their gaze towards the ground. Because this downward shift moves the stimulus image 355 upwards in the visual field, it is likely correlated with a consistent signal in the upper visual field. 356 Another possible explanation is that traffic lights are located in the upper right part of the screen, and 357 the hazard detection module is responsible for detecting red traffic lights that would require braking. 358 Together, these results suggest that the visual components of a DNN in an active driving task learn 359 similar representations to human visual processing during driving.

360 361

362

4.2 **BEV** PERCEPTION

The BEV perception module explains brain activity in punctate locations around the sensory and 363 motor cortex and the anterior IPS (Fig.3a). (Because of the punctate nature of these regions and 364 anatomical variability across subjects, they do not project to the same locations on the group-level 365 map in Fig.2b). The performance of this module in explaining somatosensory activity is likely due 366 to correlations between the layout of the local environment (e.g. an upcoming turn) and the driver's 367 control actions. Because the IPS is implicated in coordinate transformations (Grefkes et al., 2004), 368 and the BEV perception module operates in a top-down reference frame instead of an egocentric 369 one, it is likely that these two processes are correlated, enabling the BEV perception module to 370 predict IPS activity.

371 372

373

4.3 PLANNING

The planning module best explains activity in voxels in the somatosensory and motor cortices, the intraparietal sulcus (IPS), and anterior visual areas such as the extrastriate body area (EBA). Wellexplained regions of the somatosensory and motor cortices include the primary motor (M1H/F) and primary sensory (S1H/F) cortex for hand and foot, as well as supplementary hand and foot motor areas and the supplementary eye fields. These regions are directly responsible for planning



Figure 3: Features from modules implicated in social navigation explain variance in known 400 functional regions in the brain that perform similar tasks. In the LAV architecture, the bird's-401 eye-view (BEV) perception module (left) predicts the layout of the road and other vehicles around 402 the driver from a top-down viewpoint, while the prediction module (right) predicts the trajectories of other nearby vehicles. Unlike the other LAV modules, these are not directly connected to ei-403 ther sensory inputs or control outputs of the human driver, and their features are thus more abstract 404 sets of driving-relevant information relevant to deciding how to navigate around other vehicles. a): 405 In all three subjects, the BEV perception features explain variance in anterior intraparietal sulcus 406 (IPS), which is implicated in coordinate transformations. They also explain variance in punctate 407 somatomotor cortex locations. b): the prediction module explains variance in human middle tem-408 poral cortex (hMT+)/extrastriate body area (EBA) complex, and in the occipital face area (OFA) 409 and fusiform face area (FFA). The FFA and OFA are canonically implicated in the perception of 410 faces, while the EBA/hMT complex is canonically implicated in the perception of body parts and 411 biological motion. These results suggest that the LAV modules are specialized in a similar manner 412 to the functional brain regions mediating driving.

- 413
- 414
- 415 416

and producing motor actions to physically control the simulated vehicle. The IPS is involved in
 perceptual-motor coordination (Grefkes et al., 2004; Grefkes & Fink, 2005). In addition to these
 known functional regions, the planning module also explains some activity in the lateral prefrontal
 cortex.

The planning module processes local features around the subject's vehicle to output a planned tra-421 jectory based on these features. Since the driver produces control actions to follow their intended 422 trajectory, it is intuitive that the planning and control modules provide the best fit to motor and 423 supplementary motor areas — no LAV module contains explicit information about steering and ac-424 celeration outputs, and the features from the other modules are not as well-correlated with motor 425 control as the planned trajectory. The IPS is known to be responsible for coordinate transformations 426 (Grefkes et al., 2004); since like the BEV perception module the planning module operates in a 427 bird's-eye-view reference frame, it may make use of representations similar to those found in the 428 IPS. However, we note that overall activity in the IPS is better explained by the planning module 429 than the BEV perception module. It is likely that the IPS simultaneously supports multiple processes for active driving, and features in our model related to coordinate frame transforms and perceptual-430 motor integration may both be associated with activity in the IPS. Finally, the explained variance in 431 high-level visual areas is likely due to the correlations between certain semantic visual features (e.g. the presence of an intersection or other vehicles) and the local BEV features and appropriate driving plan.

434 435

436

4.4 TRAJECTORY PREDICTION

437 The trajectory prediction module explains variance in multiple anterior visual regions, including 438 the human middle temporal cortex (hMT+), extrastriate body area (EBA), occipital face area (OFA), and fusiform face area (FFA) (Fig.3b). (Similar to the well-predicted regions by the BEV perception 439 440 module, some of these regions, particularly EBA, do not project well to a group-level space). The EBA is canonically implicated in the perception of body parts (Downing et al., 2001) and biological 441 motion (Astafiev et al., 2004), the FFA in the perception of faces (Kanwisher & Yovel, 2006), and 442 hMT+ in the perception of optic flow (Morrone et al., 2000). Note that while the FFA is canonically 443 a face perception area, some evidence suggests that it may also function as a general expertise area 444 for object perception (Gauthier et al., 2000; Tarr & Gauthier, 2000). Because vehicles are an object 445 category of particular importance during driving, the FFA (and associated regions for the perception 446 of other agents) may be recruited in this capacity while driving (Ross et al., 2018). 447

Because this module is downstream of other modules that receive visual inputs, explained activity 448 in high-level vision areas may indicate alignment between the intermediate representations of other 449 agents in the LAV architecture and the representations of human drivers. However, even in its 450 best-performing regions, the trajectory prediction module explains less variance than the planning 451 module, even though the planning module does not explicitly represent information about other 452 drivers. This suggests that the representations used by the trajectory planning module may not be 453 well-aligned with those used by the human brain. Reasoning about the behavior of other agents is 454 a problem that has proven especially challenging in autonomous driving (Peters et al., 2024). Our 455 results suggest that the strategy of the LAV algorithm, which projects all agents into a common representation and predicts trajectories for the ego vehicle and other vehicles in the same way, may 456 457 not be a good match for how humans represent and reason about other drivers.

4.5 CONTROL

460 The control module best explains activity in a similar network of brain regions as the planning mod-461 ule, namely the somatosensory and motor cortices, the IPS, and some anterior visual regions. This 462 similarity is likely because the trajectory plan that the control module attempts to follow is produced 463 by the planning module and is therefore highly correlated with some of the features in the planning 464 module. In the anterior visual cortex, the controller module slightly outperforms the planning mod-465 ule in the retrosplenial cortex (RSC) and the parahippocampal place area (PPA), two visual scene 466 perception regions. This is likely because while the control and planning module features are both 467 correlated with the contents of the local scene, the trajectory plan and braking probability features 468 from the control module are more succinct than the high-dimensional representation in the plan-469 ning module, and might therefore provide a stronger correlation with certain aspects of visual scene perception encoded in the RSC and PPA. 470

471 472

473

458 459

4.6 LIMITATIONS AND FUTURE WORK

In our experimental paradigm, the human subject controls the vehicle and therefore the stimulus 474 (rendered RGB images from the driving simulator). This interactivity produces a naturalistic stimu-475 lus distribution, whereas previous work on comparing visual representations between the brain and 476 CNNs (Agrawal et al., 2014; Yamins et al., 2014) have all used tightly controlled stimuli. While 477 having an interactive task allowed us to gain insights about an active sensorimotor task, it also re-478 sults in strong correlations between task state, brain activity, and driving DNN activations. These 479 strong correlations allow for brain activity to be predicted from driving DNN activations even if 480 algorithms and representations are not aligned between the DNN and the brain. For example, if 481 the subject always stops their car when another vehicle is directly in front, then activity in the feet 482 motor areas (and therefore the accelerator/brake pedals) will be predictable from a representation that encodes the presence of another vehicle even though the feet motor areas do not actually encode 483 this information. An example of this effect is shown in Fig.?? in appendix A.2. We mitigate this 484 issue by fitting correlated models simultaneously with banded ridge regression, and then partition-485 ing the total explained variance across models (Nunez-Elizalde et al., 2019; Dupré la Tour et al.,

2022) to identify the best-performing feature space out of many correlated feature spaces that could independently explain similar amounts of variance.

These results do not tell us about the performance of our features relative to other possible feature 489 spaces that we have not included in the regression. Since our encoding model, based on features 490 from the LAV driving DNN, is only one hypothesis about the representations the brain may use for 491 driving, in future work we plan to compare encoding models based on other algorithms for driving. 492 It would be especially interesting to investigate encoding models that incorporate emerging trends 493 in deep learning for autonomous driving such as world modeling and uncertainty modeling (Chen 494 et al., 2024). Building encoding models using features from alternative driving DNNs will enable 495 us to explore which DNN network architectures and training objectives better reflect the algorithms 496 and representations used by the human brain for driving.

497 To our knowledge, we have presented the first study that compares brain activity and DNN activa-498 tions in an interactive, closed-loop, and naturalistic task. Driving requires the agent to continuously 499 perceive the stimulus, make decisions, and produce actions that in turn affects the perceived stim-500 ulus. Previous work relating DNNs to the brain have all used perceptual tasks in which the agent 501 could not interact with the stimulus, such as viewing photographs (Yamins et al., 2014), listening to 502 words or music clips (Kell et al., 2018), or categorizing an image (Flesch et al., 2022). Such tasks 503 engage a single sensory system of the brain and reflect only a very narrow subset of natural behavior. Here, we related not only perceptual representations between a neural network and the brain, 504 but also representations for using these perceptual inputs to make plans and produce action outputs. 505

5 CONCLUSION

507 508

529

538

506

509 In this work, we have presented an in-depth comparison of the cognitive representations for driving 510 in humans and in a deep end-to-end neural network for autonomous driving. By directly comparing 511 the representations in a deep learning driving algorithm with humans for the first time, we directly 512 evaluate the alignment between the human brain and the components of a driving network. Our 513 results highlight a striking similarity in the way that driving-relevant representations are organized between humans and the deep neural network. However, we also identified aspects of the network 514 that might be less aligned with humans, including the representation of other vehicles on the road. 515 While driving algorithms do not necessarily need to mimic humans to perform well, humans are still 516 more capable than autonomous vehicles in most environments. Thus, better alignment with human 517 representations may improve autonomous vehicle performance. Furthermore, better alignment with 518 human representations may enable autonomous vehicles to make better inferences about other (hu-519 man) drivers, an important component of safe and socially acceptable driving. In turn, deep neural 520 networks can provide novel insights about how the brain represents information for solving complex 521 tasks in dynamic environments. Our approach provides a promising framework for understanding 522 how both humans and AI agents solve dynamic, active tasks, both in driving as well as in other 523 domains. 524

525 REFERENCES

- 527 Dimitris Achlioptas. Database-friendly random projections: Johnson-lindenstrauss with binary
 528 coins. *Journal of computer and System Sciences*, 66(4):671–687, 2003.
- Pulkit Agrawal, Dustin Stansbury, Jitendra Malik, and Jack L Gallant. Pixels to voxels: modeling
 visual representation in the human brain. *arXiv preprint arXiv:1407.5104*, 2014.

Serguei V Astafiev, Christine M Stanley, Gordon L Shulman, and Maurizio Corbetta. Extrastriate body area in human occipital cortex responds to the performance of motor actions. *Nature neuroscience*, 7(5):542–548, 2004.

- Kaveh Azadeh, René De Koster, and Debjit Roy. Robotized and automated warehouse systems:
 Review and recent developments. *Transportation Science*, 53(4):917–945, 2019.
- Francis Bach. Sharp analysis of low-rank kernel matrix approximations. In *Conference on learning theory*, pp. 185–209. PMLR, 2013.
 - 10

- 540 Yashar Behzadi, Khaled Restom, Joy Liau, and Thomas T Liu. A component based noise correction 541 method (compcor) for bold and perfusion based fmri. Neuroimage, 37(1):90-101, 2007. 542 Matteo Carandini, Jonathan B Demb, Valerio Mante, David J Tolhurst, Yang Dan, Bruno A Ol-543 shausen, Jack L Gallant, and Nicole C Rust. Do we know what the early visual system does? 544 Journal of Neuroscience, 25(46):10577-10597, 2005. 546 Dian Chen and Philipp Krähenbühl. Learning from all vehicles. In Proceedings of the IEEE/CVF 547 Conference on Computer Vision and Pattern Recognition, pp. 17222–17231, 2022. 548 Li Chen, Penghao Wu, Kashyap Chitta, Bernhard Jaeger, Andreas Geiger, and Hongyang Li. End-549 to-end autonomous driving: Challenges and frontiers. IEEE Transactions on Pattern Analysis and 550 *Machine Intelligence*, 2024. 551 552 John Demuyakor. Ghana go digital agenda: The impact of zipline drone technology on digital emer-553 gency health delivery in ghana. Shanlax International Journal of Arts, Science and Humanities, 554 8(1):242-253, 2020. 555 Fatma Deniz, Anwar O Nunez-Elizalde, Alexander G Huth, and Jack L Gallant. The representation 556 of semantic information across human cerebral cortex during listening versus reading is invariant to stimulus modality. Journal of Neuroscience, 39(39):7722-7736, 2019. 558 559 Paul E Downing, Yuhong Jiang, Miles Shuman, and Nancy Kanwisher. A cortical area selective for visual processing of the human body. Science, 293(5539):2470-2473, 2001. 561 Tom Dupré la Tour, Michael Eickenberg, Anwar O Nunez-Elizalde, and Jack L Gallant. Feature-562 space selection with banded ridge regression. *NeuroImage*, 264:119728, 2022. 563 564 Yonina C Eldar and Gitta Kutyniok. Compressed sensing: theory and applications. Cambridge 565 university press, 2012. 566 Timo Flesch, Keno Juechems, Tsvetomira Dumbalska, Andrew Saxe, and Christopher Summerfield. 567 Orthogonal representations for robust context-dependent task performance in brains and neural 568 networks. Neuron, 110(7):1258-1270, 2022. 569 570 James S Gao, Alexander G Huth, Mark D Lescroart, and Jack L Gallant. Pycortex: an interactive surface visualizer for fmri. Frontiers in neuroinformatics, 9:23, 2015. 571 572 Isabel Gauthier, Pawel Skudlarski, John C Gore, and Adam W Anderson. Expertise for cars and 573 birds recruits brain areas involved in face recognition. Nature neuroscience, 3(2):191-197, 2000. 574 575 Christian Grefkes and Gereon R Fink. The functional organization of the intraparietal sulcus in humans and monkeys. Journal of anatomy, 207(1):3–17, 2005. 576 577 Christian Grefkes, Afra Ritzl, Karl Zilles, and Gereon R Fink. Human medial intraparietal cortex 578 subserves visuomotor coordinate transformation. Neuroimage, 23(4):1494–1506, 2004. 579 580 Joachim Gross. Magnetoencephalography in cognitive neuroscience: a primer. Neuron, 104(2): 189-204, 2019. 581 582 Umut Güçlü and Marcel AJ Van Gerven. Deep neural networks reveal a gradient in the complexity of 583 neural representations across the ventral stream. Journal of Neuroscience, 35(27):10005–10014, 584 2015. 585 Milad Haghani, Michiel CJ Bliemer, Bilal Farooq, Inhi Kim, Zhibin Li, Cheol Oh, Zahra Shah-586 hoseini, and Hamish MacDougall. Applications of brain imaging methods in driving behaviour research. Accident Analysis & Prevention, 154:106093, 2021. 588 589 Kathleen A Hansen, Kendrick N Kay, and Jack L Gallant. Topographic organization in and near 590 human visual area v4. Journal of Neuroscience, 27(44):11896–11911, 2007. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-592
 - nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

603

619

630

631

635

637

594	Alexander G Huth, Shinji Nishimoto, An T Vu, and Jack L Gallant. A continuous semantic space
595	describes the representation of thousands of object and action categories across the human brain.
596	Neuron, 76(6):1210–1224, 2012.
597	

- Alexander G Huth, Wendy A De Heer, Thomas L Griffiths, Frédéric E Theunissen, and Jack L 598 Gallant. Natural speech reveals the semantic maps that tile human cerebral cortex. Nature, 532 (7600):453-458, 2016. 600
- 601 Mark Jenkinson and Stephen Smith. The role of registration in functional magnetic resonance imag-602 ing. Medical Image Registration. CRC Press, New York, pp. 183-198, 2001.
- Mark Jenkinson, Peter Bannister, Michael Brady, and Stephen Smith. Improved optimization for the 604 robust and accurate linear registration and motion correction of brain images. *Neuroimage*, 17(2): 605 825-841, 2002. 606
- 607 Nancy Kanwisher and Galit Yovel. The fusiform face area: a cortical region specialized for the perception of faces. Philosophical Transactions of the Royal Society B: Biological Sciences, 361 608 (1476):2109-2128, 2006. 609
- 610 Alexander JE Kell, Daniel LK Yamins, Erica N Shook, Sam V Norman-Haignere, and Josh H Mc-611 Dermott. A task-optimized neural network replicates human auditory behavior, predicts brain 612 responses, and reveals a cortical processing hierarchy. Neuron, 98(3):630-644, 2018. 613
- Tom Dupre la Tour, Michael Lu, Michael Eickenberg, and Jack L Gallant. A finer mapping of 614 convolutional neural network layers to the visual cortex. In SVRHM 2021 Workshop@ NeurIPS, 615 2021. 616
- 617 David Lattanzi and Gregory Miller. Review of robotic infrastructure inspection systems. Journal of 618 Infrastructure Systems, 23(3):04017004, 2017.
- Ping Li, Trevor J Hastie, and Kenneth W Church. Very sparse random projections. In Proceedings 620 of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, 621 pp. 287-296, 2006. 622
- 623 Yichao Lu, Paramveer Dhillon, Dean P Foster, and Lyle Ungar. Faster ridge regression via the sub-624 sampled randomized hadamard transform. Advances in neural information processing systems, 26, 2013. 625
- 626 Maria Concetta Morrone, M Tosetti, D Montanaro, A Fiorentini, Giovanni Cioni, and DC Burr. A 627 cortical area that responds specifically to optic flow, revealed by fmri. Nature neuroscience, 3 628 (12):1322-1328, 2000. 629
 - Thomas Naselaris, Kendrick N Kay, Shinji Nishimoto, and Jack L Gallant. Encoding and decoding in fmri. Neuroimage, 56(2):400-410, 2011.
- 632 Jordan Navarro, Emanuelle Reynaud, and François Osiurak. Neuroergonomics of car driving: A 633 critical meta-analysis of neuroimaging data on the human brain behind the wheel. Neuroscience 634 & Biobehavioral Reviews, 95:464–479, 2018.
- Shinji Nishimoto, An T Vu, Thomas Naselaris, Yuval Benjamini, Bin Yu, and Jack L Gallant. Re-636 constructing visual experiences from brain activity evoked by natural movies. Current biology, 21(19):1641-1646, 2011. 638
- 639 Anwar O Nunez-Elizalde, Alexander G Huth, and Jack L Gallant. Voxelwise encoding models with 640 non-spherical multivariate normal priors. Neuroimage, 197:482-492, 2019.
- 641 Seiji Ogawa, Tso-Ming Lee, Alan R Kay, and David W Tank. Brain magnetic resonance imaging 642 with contrast dependent on blood oxygenation. proceedings of the National Academy of Sciences, 643 87(24):9868-9872, 1990. 644
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-645 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and 646 E. Duchesnay. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 647 12:2825-2830, 2011.

- Lasse Peters, Andrea Bajcsy, Chih-Yuan Chiu, David Fridovich-Keil, Forrest Laine, Laura Ferranti, and Javier Alonso-Mora. Contingency games for multi-agent interaction. *IEEE Robotics and Automation Letters*, 2024.
- Jonathan D Power, Benjamin M Silver, Melanie R Silverman, Eliana L Ajodan, Dienke J Bos, and Rebecca M Jones. Customized head molds reduce motion during resting state fmri scans. *NeuroImage*, 189:141–149, 2019.
- John W Pratt. Dividing the indivisible: Using simple symmetry to partition variance explained.
 In *Proceedings of the second international Tampere conference in statistics*, 1987, pp. 245–260.
 Department of Mathematical Sciences, University of Tampere, 1987.
- Eduardo Romera, José M Alvarez, Luis M Bergasa, and Roberto Arroyo. ERFNet: Efficient residual factorized convnet for real-time semantic segmentation. *IEEE Transactions on Intelligent Transportation Systems*, 19(1):263–272, 2017.
- David A Ross, Benjamin J Tamber-Rosenau, Thomas J Palmeri, JieDong Zhang, Yaoda Xu, and Is abel Gauthier. High-resolution functional magnetic resonance imaging reveals configural process ing of cars in right anterior fusiform face area of car experts. *Journal of cognitive neuroscience*, 30(7):973–984, 2018.
- Tom A Schweizer, Karen Kan, Yuwen Hung, Fred Tam, Gary Naglie, and Simon J Graham. Brain activity during driving with distraction: an immersive fmri study. *Frontiers in human neuroscience*, 7:53, 2013.
- Hugo J Spiers and Eleanor A Maguire. Neural substrates of driving behaviour. *Neuroimage*, 36(1): 245–255, 2007.
- Yu Takagi and Shinji Nishimoto. High-resolution image reconstruction with latent diffusion models
 from human brain activity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14453–14463, 2023.
- Michael J Tarr and Isabel Gauthier. Ffa: a flexible fusiform area for subordinate-level visual processing automatized by expertise. *Nature neuroscience*, 3(8):764–769, 2000.
- David P Woodruff et al. Sketching as a tool for numerical linear algebra. *Foundations and Trends*® *in Theoretical Computer Science*, 10(1–2):1–157, 2014.
- Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the national academy of sciences*, 111(23):8619–8624, 2014.
 - Yu Zhang, Peter Tiňo, Aleš Leonardis, and Ke Tang. A survey on neural network interpretability. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 5(5):726–742, 2021.

699 700

696 697

651

658

669

675

680

684

685

701