

Toward Goal-Driven Neural Network Models for the Rodent Whisker-Trigeminal System

Chengxu Zhuang
Department of Psychology
Stanford University
Stanford, CA 94305
chengxuz@stanford.edu

May 31, 2017

In large part, rodents “see” the world through their whiskers, a powerful tactile sense enabled by a series of brain areas that form the whisker-trigeminal system. Raw sensory data arrives in the form of mechanical input to the exquisitely sensitive, actively-controllable whisker array, and is processed through a sequence of neural circuits, eventually arriving in cortical regions that communicate with decision-making and memory areas. Although a long history of experimental studies has characterized many aspects of these processing stages, the computational operations of the whisker-trigeminal system remain largely unknown. In the present work, we take a *goal-driven* deep neural network (DNN) approach to modeling these computations. First, we construct a biophysically-realistic model of the rat whisker array. We then generate a large dataset of whisker sweeps across a wide variety of 3D objects in highly-varying poses, angles, and speeds. Next, we train DNNs from several distinct architectural families to solve a shape recognition task in this dataset. Each architectural family represents a structurally-distinct hypothesis for processing in the whisker-trigeminal system, corresponding to different ways in which spatial and temporal information can be integrated. We find that most networks perform poorly on the challenging shape recognition task, but that specific architectures from several families can achieve reasonable performance levels. Finally, we show that Representative Dissimilarity Matrices (RDMs), a tool for comparing population codes between neural systems, can separate these higher-performing net-

works with data of a type that could plausibly be collected in a neurophysiological or imaging experiment. Our results are a proof-of-concept that DNN models of the whisker-trigeminal system are potentially within reach.

1 Introduction

The sensory systems of brains do remarkable work in extracting behaviorally useful information from noisy and complex raw sense data. Vision systems process intensities from retinal photoreceptor arrays, auditory systems interpret the amplitudes and frequencies of hair-cell displacements, and somatosensory systems integrate data from direct physical interactions. [28] Although these systems differ radically in their input modalities, total number of neurons, and specific neuronal microcircuits, they share two fundamental characteristics. First, they are hierarchical sensory cascades, albeit with extensive feedback, consisting of sequential processing stages that together produce a complex transformation of the input data. Second, they operate in inherently highly-structured spatiotemporal domains, and are generally organized in maps that reflect this structure [11].

Extensive experimental work in the rodent whisker-trigeminal system has provided insights into how these principles help rodents use their whiskers (also known as *vibrissae*) to tactually explore objects in their environment. Similar to hierarchical processing in the visual system (e.g., from V1 to V2, V4 and IT [11, 12]), processing in the somatosensory system is also known to be hierarchical[27, 17, 18]. For example, in the whisker trigeminal system, information from the whiskers is relayed from primary sensory neurons in the trigeminal ganglion to multiple trigeminal nuclei; these nuclei are the origin of several parallel pathways conveying information to the thalamus [36, 24] and then to primary and secondary somatosensory cortex (S1 and S2) [4]. However, although the rodent somatosensory system has been the subject of extensive experimental efforts[2, 26, 20, 32], there have been comparatively few attempts at computational modeling of this important sensory system.

Recent work has shown that deep neural networks (DNNs), whose architectures inherently contain hierarchy and spatial structure, can be effective models of neural processing in vision[34, 21] and audition[19]. Motivated by these successes, in this work we illustrate initial steps toward using DNNs to model rodent somatosensory systems. Our driving hypothesis is that the vibrissal-trigeminal system is optimized to use whisker-based sensor data to solve somatosensory shape-recognition tasks in complex, variable real-world environments. The underlying idea of this approach is thus to use *goal-driven* modeling (Fig 1), in which the DNN parameters — both discrete and continuous — are optimized for performance on a challenging ethologically-relevant task[35]. Insofar as shape recognition is a strong constraint on network parameters, optimized neural networks resulting from such a task may be an effective model of real trigeminal-system neural response patterns.

This idea is conceptually straightforward, but implementing it involves sur-

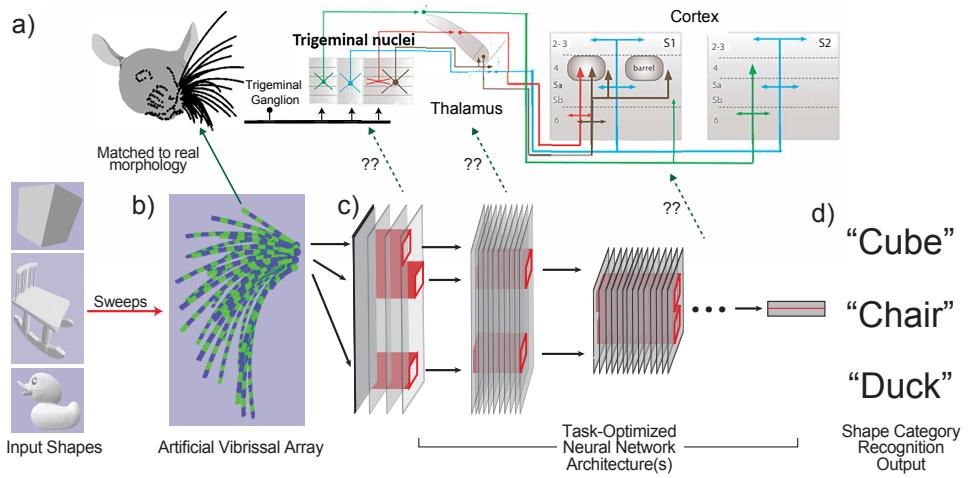


Figure 1: Goal-Driven Approach to Modeling Barrel Cortex: **a.** Rodents have highly sensitive whisker (vibrissal) arrays that provide input data about the environment. Mechanical signals from the vibrissae are relayed by primary sensory neurons of the trigeminal ganglion to the trigeminal nuclei, the origin of multiple parallel pathways to S1 and S2. (Figure modified from [8].) This system is a prime target for modeling because it is likely to be richly representational, but its computational underpinnings are largely unknown. Our long-term approach to modeling the whisker-trigeminal system is *goal-driven*: using an artificial whisker-array input device built using extensive biophysical measurements (**b.**), we seek to optimize neural networks of various architectures (**c.**) to solve ethologically-relevant shape recognition tasks (**d.**), and then measure to the extent to which these networks predict fine-grained response patterns in real neural recordings.

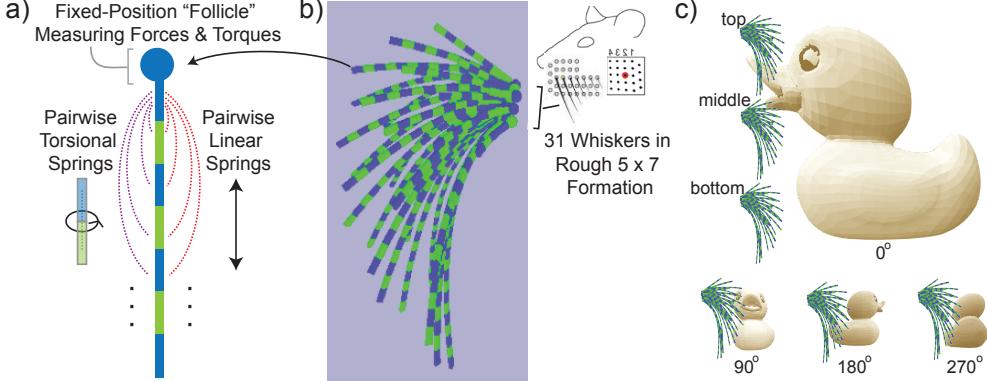


Figure 2: Dynamic Three-Dimensional Whisker Model: **a.** Each whisker element is composed of a set of cuboid links. The follicle cuboid has a fixed location, and is attached to movable cuboids making up the rest of the whisker. Motion is constrained by linear and torsional springs between each pair of cuboids. The number of cuboid links and spring equilibrium displacements are chosen to match known whisker length and curvature [31], while damping and spring stiffness parameters are chosen to ensure mechanically plausible whisker motion trajectories. **b.** We constructed a 31-whisker array, arranged in a rough 5x7 grid (with 4 missing elements) on an ellipsoid representing the rodent’s mystacial pad. Whisker number and placement was matched to the known anatomy of the rat [31]. **c.** During dataset construction, the array is brought into contact with each object at three vertical heights, and four 90°-separated angles, for a total of 12 sweeps. The object’s size, initial orientation angle, as well as sweep speed, vary randomly between each group of 12 sweeps. Forces and torques are recorded at the three cuboids closest to the follicle, for a total of 18 measurements per whisker at each timepoint.

mounting several challenges. Unlike vision or audition, where signals from the retina or cochlea can for many purposes be approximated by a simple structure (namely, a uniform data array representing light or sound intensities and frequencies), the equivalent mapping from stimulus (e.g. object in a scene) to sensor input in the whisker system is much less direct. Thus, a biophysically-realistic embodied model of the whisker array is a critical first component of any model of the vibrissal system. Once the sensor array is available, a second key problem is building a neural network that can accept whisker data input and use it to solve relevant tasks. Aside from the question of the neural network design itself, knowing what the “relevant tasks” are for training a rodent whisker system, in a way that is sufficiently concrete to be practically actionable, is a significant unknown, given the very limited amount of ethologically-relevant behavioral data on rodent sensory capacities[32, 22, 25, 1, 9]. Collecting neural data of sufficient coverage and resolution to quantitatively evaluate one or more task-optimized neural network models represents a third major challenge. In this work, we show initial steps toward the first two of these problems (sensor modeling and neural network design/training).

2 Modeling the Whisker Array Sensor

In order to provide our neural networks inputs similar to those of the rodent vibrissal system, we constructed a physically-realistic three-dimensional (3D) model of the rodent vibrissal array (Fig. 2). To help ensure biological realism,

we used an anatomical model of the rat head and whisker array that quantifies whisker number, length, and intrinsic curvature as well as relative position and orientation on the rat’s face [31]. We wanted the mechanics of each whisker to be reasonably accurate, but at the same time, also needed simulations to be fast enough to generate a large training dataset. We therefore used the Bullet [33], an open-source real-time physics engine used in many video games.

Statics. Individual whiskers were each modeled as chains of “cuboid” links with a square cross-section and length of 2mm. The number of links in each whisker was chosen to ensure that the total whisker length matched that of the corresponding real whisker (Fig. 2 a). The first (most proximal) link of each simulated whisker corresponded to the follicle at the whisker base, where the whisker inserts into the rodent’s face. Each whisker follicle was fixed to a single location in 3D space. The links of the whisker are given first-order linear and rotational damping factors to ensure that unforced motions dissipate over time. To simplify the model, the damping factors were assumed to be the same across all links of a given whisker, but different from whisker to whisker. Each pair of links within a whisker was connected with linear and torsional first-order springs; these springs both have two parameters (equilibrium displacement and stiffness). The equilibrium displacements of each spring were chosen to ensure that the whisker’s overall static shape matched the measured curvature for the corresponding real whisker. Although we did not specifically seek to match the detailed biophysics of the whisker mechanics (e.g. the fact that the stiffness of the whisker increases with the 4th power of its radius), we assumed that the stiffness of the springs spanning a given length were linearly correlated to the distance between the starting position of the spring and the base, roughly capturing the fact that the whisker is thicker and stiffer at the bottom [13].

The full simulated whisker array consisted of 31 simulated whiskers, ranging in length from 8mm to 60mm (Fig. 2b). The fixed locations of the follicles of the simulated whiskers were placed on a curved ellipsoid surface modeling the rat’s mystacial pad (cheek), with the relative locations of the follicles on this surface obtained from the morphological model [31], forming roughly a 5×7 grid-like pattern with four vacant positions.

Dynamics. Whisker dynamics are generated by collisions with moving three-dimensional rigid bodies, also modeled as Bullet physics objects. The motion of a simulated whisker in reaction to external forces from a collision is constrained only by the fixed spatial location of the follicle, and by the damped dynamics of the springs at each node of the whisker. However, although the spring equilibrium displacements are determined by static measurements as described above, the damping factors and spring stiffnesses cannot be fully determined from these data. If we had detailed dynamic trajectories for all whiskers during realistic motions (e.g. [29]), we would have used this data to determine these parameters, but such data are not yet available.

In the absence of empirical trajectories, we used a heuristic method to determine damping and stiffness parameters, maximizing the “mechanical plausibility” of whisker behavior. Specifically, we constructed a battery of scenarios in which forces were applied to each whisker for a fixed duration. These scenarios included

pushing the whisker tip towards its base (axial loading), as well as pushing the whisker parallel or perpendicular to its intrinsic curvature (transverse loading in or out of the plane of intrinsic curvature). For each scenario and each potential setting of the unknown parameters, we simulated the whisker’s recovery after the force was removed, measuring the maximum displacement between the whisker base and tip caused by the force prior to recovery (d), the total time to recovery (T), the average arc length travelled by each cuboid during recovery (S), and the average translational speed of each cuboid during recovery (v). We used metaparameter optimization [3] to automatically identify stiffness and damping parameters that simultaneously minimized the time and complexity of the recovery trajectory, while also allowing the whisker to be flexible. Specifically, we minimized the loss function $0.025S + d + 20T - 2v$, where the coefficients were set to make terms of comparable magnitude. The optimization was performed for every whisker independently, as whisker length and curvature interacts nonlinearly with its recovery dynamics.

3 A Large-Scale Whisker Sweep Dataset

Using the whisker array, we generated a dataset of whisker responses to a variety of objects.

Sweep Configuration. The dataset consists of series of simulated sweeps, mimicking one action in which the rat runs its whiskers past an object while holding its whiskers fixed (no active whisking). During each sweep, a single 3D object moves through the whisker array from front to back (rostral to caudal) at a constant speed. Each sweep lasts a total of one second, and data is sampled at 110Hz. Sweep scenarios vary both in terms of the identity of the object presented, as well as the position, angle, scale (defined as the length of longest axis), and speed at which it is presented. To simulate observed rat whisking behavior in which animals often sample an object at several vertical locations (head pitches) [14], sweeps are performed at three different heights along the vertical axis and at each of four positions around the object (0° , 90° , 180° , and 270° around the vertical axis), for a total of 12 sweeps per object/latent variable setting (Fig. 2c).

Latent variables settings are sampled randomly and independently on each group of sweeps, with object rotation sampled uniformly within the space of all 3D rotations, object scale sampled uniformly between 25-135mm, and sweep speed sampled randomly between 77-154mm/s. Once these variables are chosen, the object is placed at a position that is chosen uniformly in a $20 \times 8 \times 20\text{mm}^3$ volume centered in front of the whisker array at the chosen vertical height, and is moved along the ray toward the center of the whisker array at the chosen speed. The position of the object may be adjusted to avoid collisions with the fixed whisker base ellipsoid during the sweep. See implementary information for details.

The data collected during a sweep includes, for each whisker, the forces and torques from all springs connecting to the three cuboids most proximate to the base of the whisker. This choice reflects the idea that mechanoreceptors are distributed along the entire length of the follicle at the whisker base [10]. The

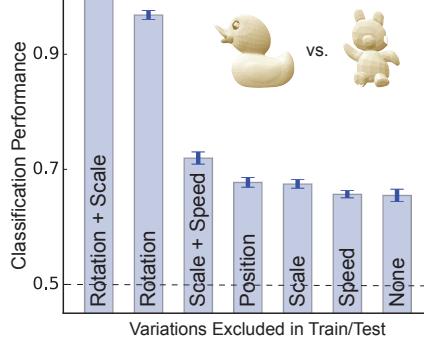


Figure 3: Basic validating experiment: Basic validation of performance of binary linear classifier trained on raw sensor output to distinguish between two shapes (in this case, a duck versus a teddy bear). The classifier was trained/tested on several equal-sized datasets in which variation on one or more latent variable axes is suppressed. “None” indicates that all variations are present. Dotted line represents chance performance (50%).

collected data comprises a matrix of shape $110 \times 31 \times 3 \times 2 \times 3$, with dimensions respectively corresponding to: the 110 time samples; the 31 spatially distinct whiskers; the 3 recorded cuboids; the forces and torques from each cuboid; and the three directional components of force/torque.

Object Set. The objects used in each sweep are chosen from a subset of the ShapeNet [6] dataset, which contains over 50,000 3D objects, each with a distinct geometry, belonging to 55 categories. Because the 55 ShapeNet categories are at a variety of levels of within-category semantic similarity, we refined the original 55 categories into a taxonomy of 117 (sub)categories that we felt had a more uniform amount of within-category shape similarity. The distribution of number of ShapeNet objects is highly non-uniform across categories, so we randomly subsampled objects from large categories. This procedure ensured that all categories contained approximately the same number of objects. Our final object set included 9,981 objects in 117 categories, ranging between 41 and 91 object exemplars per category (mean=85.3, median=91, std=10.2, see supplementary material for more details). To create the final dataset, for every object, 26 independent samples of rotation, scaling, and speed were drawn and the corresponding group of 12 sweeps created. Out of these 26 sweep groups, 24 were added to a training subset, while the remainder were reserved for testing.

Basic Sensor Validation. To confirm that the whisker array was minimally functional before proceeding to more complex models, we produced smaller versions of our dataset in which sweeps were sampled densely for two objects (a bear and a duck). We also produced multiple easier versions of this dataset in which variation along one or several latent variables was suppressed. We then trained binary support vector machine (SVM) classifiers to report object identity in these datasets, using only the raw sensor data as input, and testing classification accuracy on held-out sweeps (Fig. 3). We found that with scale and object rotation variability suppressed (but with speed and position variability retained), the sensor was able to nearly perfectly identify the objects. However,

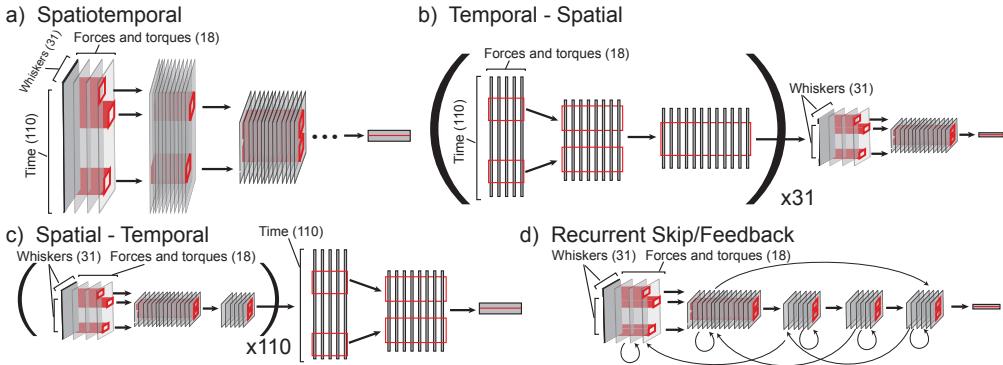


Figure 4: Families of DNN Architectures tested: **a.** “Spatiotemporal” models include spatiotemporal integration at all stages. Convolution is performed on both spatial and temporal data dimensions, followed by one or several fully connected layers. **b.** “Temporal-Spatial” networks in which temporal integration is performed separately before spatial integration. Temporal integration consists of one-dimensional convolution over the temporal dimension, separately for each whisker. In spatial integration stages, outputs from each whisker are registered to their natural two-dimensional (2D) spatial grid and spatial convolution performed. **c.** In “Spatial-Temporal” networks, spatial convolution is performed first, replicated with shared weights across time points; this is then followed by temporal convolution. **d.** Recurrent networks do not explicitly contain separate units to handle different discrete timepoints, relying instead on the states of the units to encode memory traces. These networks can have local recurrence (e.g. simple addition or more complicated motifs like LSTMs or GRUs), as well as long-range skip and feedback connections.

with all sources of variability present, the SVM was just above chance in its performance, while combinations of variability are more challenging for the sensor than others (details can be found in supplementary information). Thus, we concluded that our virtual whisker array was basically functional, but that unprocessed sensor data cannot be used to directly read out object shape in anything but the most highly controlled circumstances. As in the case of vision, it is exactly this circumstance that calls for a deep cascade of sensory processing stages.

4 Computational Architectures

We trained deep neural networks (DNNs) in a variety of different architectural families (Fig. 4). These architectural families represent qualitatively different classes of hypotheses about the computations performed by the stages of processing in the vibrissal-trigeminal system. The fundamental questions explored by these hypotheses are how and where temporal and spatial information are integrated. Within each architectural family, the differences between specific parameter settings represent nuanced refinements of the larger hypothesis of that family. Parameter specifics include how many layers of each type are in the network, how many units are allocated to each layer, what kernel sizes are used at each layer, and so on. Biologically, these parameters may correspond to the number of brain regions (areas) involved, how many neurons these regions have relative to each other, and neurons’ local spatiotemporal receptive field sizes [35].

Simultaneous Spatiotemporal Integration. In this family of networks (Fig. 4a), networks consisted of convolution layers followed by one or more fully connected layers. Convolution is performed simultaneously on both temporal and spatial dimensions of the input (and their corresponding downstream dimensions). In other words, temporally-proximal responses from spatially-proximal whiskers are combined together simultaneously, so that neurons in each successive layers have larger receptive fields in both spatial and temporal dimensions at once. We evaluated both 2D convolution, in which the spatial dimension is indexed linearly across the list of whiskers (first by vertical columns and then by lateral row on the 5×7 grid), as well as 3D convolution in which the two dimensions of the 5×7 spatial grid are explicitly represented. Data from the three vertical sweeps of the same object were then combined to produce the final output, culminating in a standard softmax cross-entropy.

Separate Spatial and Temporal Integration. In these families, networks begin by integrating temporal and spatial information separately (Fig. 4b-c). One subclass of these networks are “Temporal-Spatial” (Fig. 4b), which first integrate temporal information for each individual whisker separately and then combine the information from different whiskers in higher layers. Temporal processing is implemented as 1-dimensional convolution over the temporal dimension. After several layers of temporal-only processing (the number of which is a parameter), the outputs at each whisker are then reshaped into vectors and combined into a 5×7 whisker grid. Spatial convolutions are then applied for several layers. Finally, as with the spatiotemporal network described above, features from three sweeps are concatenated into a single fully connected layer which outputs softmax logits.

Conversely, “Spatial-Temporal” networks (Fig. 4c) first use 2D convolution to integrate across whiskers for some number of layers, with shared parameters between the copies of the network for each timepoint. The temporal sequence of outputs is then combined, and several layers of 1D convolution are then applied in the temporal domain. Both Temporal-Spatial and Spatial-Temporal networks can be viewed as subclasses of 3D simultaneous spatiotemporal integration in which initial and final portions of the network have kernel size 1 in the relevant dimensions. These two network families can thus be thought of as two different strategies for allocating parameters between dimensions, i.e. different possible biological circuit structures.

Recurrent Neural Networks with Skip and Feedback Connections. This family of networks (Fig. 4d) does not allocate units or parameters explicitly for the temporal dimension, and instead requires temporal processing to occur via the temporal update evolution of the system. These networks are built around a core feedforward 2D spatial convolution structure, with the addition of (i) local recurrent connections, (ii) long-range feedforward skips between non-neighboring layers, and (iii) long-range feedback connections. The most basic update rule for the dynamic trajectory of such a network through (discrete) time is: $H_{t+1}^i = F_i \left(\bigoplus_{j \neq i} R_t^j \right) + \tau_i H_t^i$ and $R_t^i = A_i[H_t^i]$, where R_t^i and H_t^i are the output and hidden state of layer i at time t respectively, τ_i are decay constants, \oplus

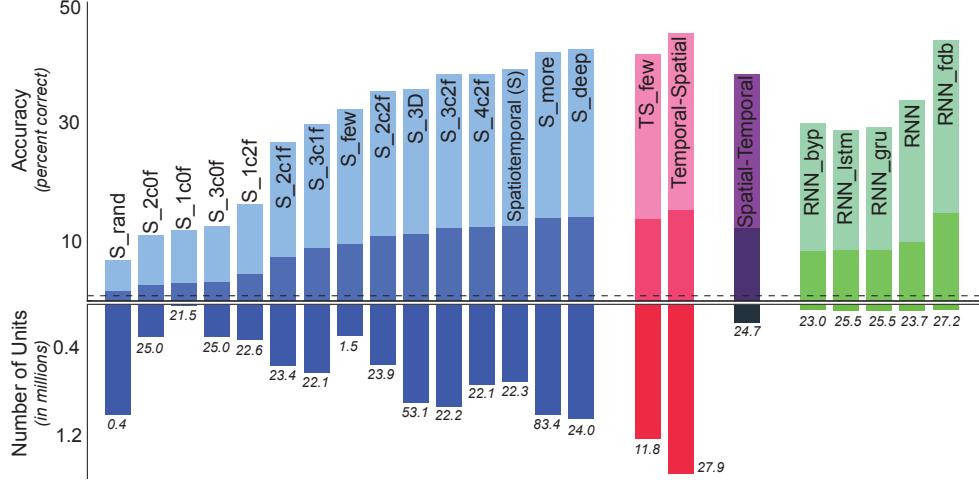


Figure 5: Performance results. Each bar in this figure represents one model. The positive y -axis is performance measured in percent correct (top1=dark bar, chance=0.85%, top5=light bar, chance=4.2%). The negative y -axis indicates the number of units in networks, in millions of units. Small italic numbers indicate number of model parameters, in millions. Model architecture family is indicated by color. "ncmf" means n convolution and m fully connected layers. Detailed definition of individual model labels can be found in supplementary material.

represents concatenation across the channel dimension with appropriate resizing to align dimensions, F_i is the standard neural network update function (e.g. 2-D convolution), and A_i is activation function at layer i . The learned parameters of this type of network include the values of the parameters of F_i , which comprises both the feedforward and feedback weights from connections coming in to layer i , as well as the decay constants τ_i . More sophisticated dynamics can be incorporated by replacing the simple additive rule above with a local recurrent structure such as Long Short-Term Memory (LSTM) [15] or Gated Recurrent Networks (GRUs) [7].

5 Results

Model Performance: Our strategy in identifying potential models of the whisker-trigeminal system is to explore many specific architectures within each architecture family, evaluating each specific architecture both in terms of its ability to solve the shape recognition task in our training dataset, and its efficiency (number of parameters and number of overall units). Because it can be misleading to compare models with different numbers of parameters, we generally evaluated models with similar numbers of parameters: exceptions are noted where they occur. As we evaluated many individual structures within each family, a list of the specific models and parameters are given in the supplementary materials.

Our results (Fig. 5) can be summarized with following conclusions:

- Many specific network choices within all families do a poor job at the task, achieving just-above-chance performance.
- However, within each family, certain specific choices of parameters lead to

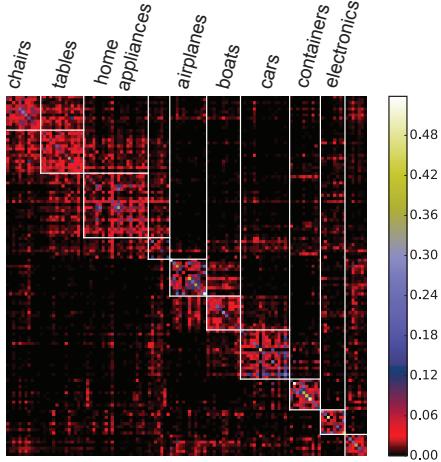


Figure 6: Confusion Matrix. Confusion Matrix for the highest-performing model (in the Temporal-Spatial family). The objects are regrouped using methods described in supplementary material.

much better network performance. Overall, the best performance was obtained for the Temporal-Spatial model, with 15.2% top-1 and 44.8% top-5 accuracy. Visualizing a confusion matrix for this network (Fig. 6) and other high-performing networks indicate that the errors they make are generally reasonable.

- Training the filters was extremely important for performance; no architecture with random filters performed above chance levels.
- Architecture depth was an important factor in performance. Architectures with fewer than four layers achieved substantially lower performance than somewhat deeper ones.
- Number of model parameters was a somewhat important factor in performance within an architectural family, but only to a point, and not between architectural families. The Temporal-Spatial architecture was able to outperform other classes while using significantly fewer parameters.
- Recurrent networks with long-range feedback were able to perform nearly as well as the Temporal-Spatial model with equivalent numbers of parameters, while using far fewer units. These long-range feedbacks appeared critical to performance, with purely local recurrent architectures (including LSTM and GRU) achieving significantly worse results.

Model Discrimination: The above results indicated that we had identified several high-performing networks in quite distinct architecture families. In other words, the strong performance constraint allows us to identify several specific candidate model networks for the biological system, reducing a much larger set of mostly non-performing neural networks into a “shortlist”. The key biologically relevant follow-up question is then: how should we distinguish between the elements in the shortlist? That is, what signatures of the differences between

these architectures could be extracted from data obtainable from experiments that use today’s neurophysiological tools?

To address this question, we used Representational Dissimilarity Matrix (RDM) analysis [23]. For a set of stimuli S , RDMs are $|S| \times |S|$ -shaped correlation distance matrices taken over the feature dimensions of a representation, e.g. matrices with ij -th entry $RDM[i, j] = 1 - \text{corr}(F[i], F[j])$ for stimuli i, j and corresponding feature output $F[i], F[j]$. The RDM characterizes the geometry of stimulus representation in a way that is independent of the individual feature dimensions. RDMs can thus be quantitatively compared between different feature representations of the same data. RDMs have been useful in establishing connections between deep neural networks and neural data from the ventral visual stream [5, 34, 21]. RDMs are readily computable from neural response pattern data samples, and are in general are comparatively robust to variability due to experimental randomness (e.g. electrode/voxel sampling). We obtained RDMs for several of our high-performing models, computing RDMs separately for each model layer (Fig. 7), averaging feature vectors over different sweeps of the same object before computing the correlations. This procedure lead to 9981×9981 -sized matrices (there were 9,981 distinct object in our dataset). We then computed distances between each layer of each model in RDM space, as in (e.g.) [21]. To determine if differences in this space between models and/or layers were significant, we computed RDMs for multiple instances of each model trained with different initial conditions, and compared the between-model to within-model distances. We found that while the top layers of models partially converged (likely because they were all trained on the same task), intermediate layers diverged substantially between models, by amounts larger than either the initial-condition-induced variability within a model layer or the distance between nearby layers of the same model (Fig. 8). This result is not entirely surprising since the models are quite structurally distinct. However, it establishes the initial computational validity for a conceptually simple and potentially feasible neuroscience experiment to distinguish between models: RDMs computed from neural recordings in multiple areas of the whisker-trigeminal system could be compared to model RDMs.

6 Conclusion

We have introduced a model of the rodent whisker array informed by biophysical data, and used it to generate a large high-variability synthetic sweep dataset. While the raw sensor data is sufficiently powerful to separate objects at low amounts of variability, at higher variation levels deeper non-linear neural networks are required to extract object identity. We found further that while many particular network architectures, especially shallow ones, fail to solve the shape recognition task, reasonable performance levels can be obtained for specific architectures within each distinct network structural family tested. We then showed that a population-level measurement that is in principle experimentally obtainable can distinguish between these higher-performing networks. To summarize, we have shown that a goal-driven DNN approach to modeling the whisker-trigeminal system is feasible. Code for all results, including the whisker model and neural

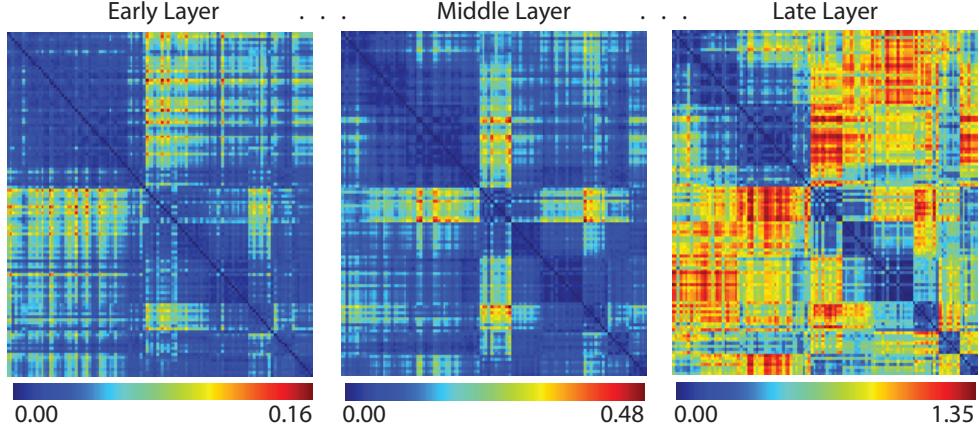


Figure 7: **RDMs visualization for High-Performing Models.** Representational Dissimilarity Matrices (RDMs) for selected layers of a high-performing network from Fig. 5, showing early, intermediate and late model layers. Model feature vectors are averaged over classes in the dataset prior to RDM computation, and RDMs are shown using the same ordering as in Fig. 6.

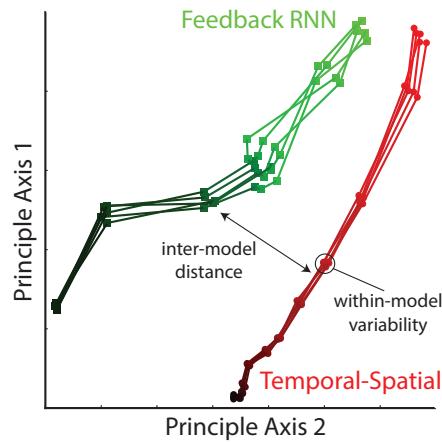


Figure 8: **Using RDMs to Discriminate Between High-Performing Models.** Two-dimensional MDS embedding of RDMs for the feedback RNN (green squares) and Temporal-Spatial (red circles) model. Points correspond to layers, lines are drawn between adjacent layers, with darker color indicating earlier layers. Multiple lines are models trained from different initial conditions, allowing within-model noise estimate.

networks, is publicly available at <https://github.com/neuroailab/barrel>.

We emphasize that the present work is proof-of-concept rather than a model of the real nervous system. A number of critical issues must be overcome before our true goal — a full integration of modeling with experimental data — becomes possible. First, although our sensor model was biophysically informed, it does not include active whisking, and the mechanical signals at the whisker bases are approximate [29, 16].

An equally important problem is that the goal that we set for our network, i.e. shape discrimination between 117 human-recognizable object classes, is not directly ethologically relevant to rodents. The primary reason for this task choice was practical: ShapeNet is a readily available and high-variability source of 3D objects. If we had instead used a small, manually constructed, set of highly simplified objects that we hoped were more “rat-relevant”, it is likely that our task would have been too simple to constrain neural networks at the scale of the real whisker-trigeminal system. Extrapolating from modeling of the visual system, training a deep net on 1000 image categories yields a feature basis that can readily distinguish between previously-unobserved categories [34, 5, 30]. Similarly, we suggest that the large and variable object set used here may provide a meaningful constraint on network structure, as the specific object geometries may be less important than having a wide spectrum of such geometries. However, a key next priority is systematically building an appropriately large and variable set of objects, textures or other class boundaries that more realistically model the tasks that a rodent faces. The specific results obtained (e.g. which families are better than others, and the exact structure of learned representations) are likely to change significantly when these improvements are made.

In the longer term, we expect to use detailed encoding models of the whisker-trigeminal system as a platform for investigating issues of representation learning and sensory-based decision making in the rodent. A particularly attractive option is to go beyond fixed class discrimination problems and situate a synthetic whisker system on a mobile animal in a navigational environment where it will be faced with a variety of actively-controlled discrete and continuous estimation problems. By doing this work with a rich sensory domain in rodents, we seek to leverage the sophisticated neuroscience tools available in these systems to go beyond what might be possible in other model systems.

7 Acknowledgement

Thanks to my advisor Daniel Yamins for tremendous help and suggestions. Also, this project is a joint work with Jonas Kubilius from Department of Brain and Cognitive Sciences of MIT and Mitra Hartmann from Departments of Biomedical Engineering and Mechanical Engineering of Northwestern University.

References

- [1] Ehsan Arabzadeh, Erik Zorzin, and Mathew E. Diamond. Neuronal encoding of texture in the whisker sensory pathway. *PLoS Biology*, 3(1), 2005.

- [2] Michael Armstrong-James, KEVIN Fox, and Ashis Das-Gupta. Flow of excitation within rat barrel cortex on striking a single vibrissa. *Journal of neurophysiology*, 68(4):1345–1358, 1992.
- [3] James Bergstra, Dan Yamins, and David D Cox. Hyperopt: A python library for optimizing the hyperparameters of machine learning algorithms. In *Proceedings of the 12th Python in Science Conference*, pages 13–20. Citeseer, 2013.
- [4] Laurens WJ Bosman, Arthur R Houweling, Cullen B Owens, Nouk Tanke, Olesya T Shevchouk, Negah Rahmati, Wouter HT Teunissen, Chiheng Ju, Wei Gong, Sebastiaan KE Koekkoek, et al. Anatomical pathways involved in generating and sensing rhythmic whisker movements. *Frontiers in integrative neuroscience*, 5:53, 2011.
- [5] Charles F Cadieu, Ha Hong, Daniel LK Yamins, Nicolas Pinto, Diego Ardila, Ethan A Solomon, Najib J Majaj, and James J DiCarlo. Deep neural networks rival the representation of primate cortex for core visual object recognition. *PLoS computational biology*, 10(12):e1003963, 2014.
- [6] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. *ArXiv*, 2015.
- [7] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
- [8] Martin Deschenes and Nadia Urbain. Vibrissal afferents from trigeminal to cortices. *Scholarpedia*, 4(5):7454, 2009.
- [9] Mathew E Diamond, Moritz von Heimendahl, Per Magne Knutson, David Kleinfeld, and Ehud Ahissar. ‘Where’ and ‘what’ in the whisker sensorimotor system. *Nat Rev Neurosci*, 9(8):601–612, 2008.
- [10] Satomi Ebara, Kenzo Kumamoto, Tadao Matsuura, Joseph E Mazurkiewicz, and Frank L Rice. Similarities and differences in the innervation of mystacial vibrissal follicle–sinus complexes in the rat and cat: a confocal microscopic study. *Journal of Comparative Neurology*, 449(2):103–119, 2002.
- [11] Daniel J Felleman and David C Van Essen. Distributed hierarchical processing in the primate cerebral cortex. *Cerebral cortex*, 1(1):1–47, 1991.
- [12] Melvyn A. Goodale and A. David Milner. Separate visual pathways for perception and action. *Trends in Neurosciences*, 15(1):20–25, 1992.
- [13] M. Hartmann. Vibrissa mechanical properties. *Scholarpedia*, 10(5):6636, 2015. revision #151934.
- [14] Jennifer A Hobbs, R Blythe Towal, and Mitra JZ Hartmann. Spatiotemporal patterns of contact across the rat vibrissal array during exploratory behavior. *Frontiers in behavioral neuroscience*, 9, 2015.
- [15] Sepp Hochreiter and Jurgen Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1–32, 1997.
- [16] Lucie A. Huet and Mitra J Z Hartmann. Simulations of a Vibrissa Slipping along a Straight Edge and an Analysis of Frictional Effects during Whisking. *IEEE Transactions on Haptics*, 9(2):158–169, 2016.

- [17] Koji Inui, Xiaohong Wang, Yohei Tamura, Yoshiki Kaneoke, and Ryusuke Kakigi. Serial processing in the human somatosensory system. *Cerebral Cortex*, 14(8):851–857, 2004.
- [18] Yoshiaki Iwamura. Hierarchical somatosensory processing. *Current Opinion in Neurobiology*, 8(4):522–528, 1998.
- [19] A *Kell, D *Yamins, S Norman-Haignere, and J McDermott. Functional organization of auditory cortex revealed by neural networks optimized for auditory tasks. In *Society for Neuroscience*, 2015.
- [20] Jason ND Kerr, Christiaan PJ De Kock, David S Greenberg, Randy M Bruno, Bert Sakmann, and Fritjof Helmchen. Spatial organization of neuronal population responses in layer 2/3 of rat barrel cortex. *Journal of neuroscience*, 27(48):13316–13328, 2007.
- [21] Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but not unsupervised, models may explain it cortical representation. *PLoS Comput Biol*, 10(11):e1003915, 2014.
- [22] Per Magne Knutsen, Maciej Pietr, and Ehud Ahissar. Haptic object localization in the vibrissal system: behavior and performance. *The Journal of neuroscience : the official journal of the Society for Neuroscience*, 26(33):8451–64, 2006.
- [23] Nikolaus Kriegeskorte, Marieke Mur, and Peter a. Bandettini. Representational similarity analysis - connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2(November):4, 2008.
- [24] Jeffrey D Moore, Nicole Mercer Lindsay, Martin Deschênes, and David Kleinfeld. Vibrissa self-motion and touch are reliably encoded along the same somatosensory pathway from brainstem through thalamus. *PLoS Biol*, 13(9):e1002253, 2015.
- [25] Daniel H. O'Connor, Simon P. Peron, Daniel Huber, and Karel Svoboda. Neural activity in barrel cortex underlying vibrissa-based object localization in mice. *Neuron*, 67(6):1048–1061, 2010.
- [26] Carl CH Petersen, Amiram Grinvald, and Bert Sakmann. Spatiotemporal dynamics of sensory responses in layer 2/3 of rat barrel cortex measured in vivo by voltage-sensitive dye imaging combined with whole-cell voltage recordings and neuron reconstructions. *Journal of neuroscience*, 23(4):1298–1309, 2003.
- [27] T P Pons, P E Garraghty, David P Friedman, and Mortimer Mishkin. Physiological evidence for serial processing in somatosensory cortex. *Science (New York, N.Y.)*, 237(4813):417–420, 1987.
- [28] Dale Purves, George J Augustine, David Fitzpatrick, Lawrence C Katz, Anthony-Samuel LaMantia, James O McNamara, and S Mark Williams. Neuroscience. *Sunderland, MA: Sinauer Associates*, 3, 2001.
- [29] Brian W Quist, Vlad Seghete, Lucie A Huet, Todd D Murphrey, and Mitra J Z Hartmann. Modeling Forces and Moments at the Base of a Rat Vibrissa during Noncontact Whisking and Whisking against an Object. *J Neurosci*, 34(30):9828–9844, 2014.
- [30] Ali S Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2014 IEEE Conference on*, pages 512–519. IEEE, 2014.

- [31] R. Blythe Towal, Brian W. Quist, Venkatesh Gopal, Joseph H. Solomon, and Mitra J Z Hartmann. The morphology of the rat vibrissal array: A model for quantifying spatiotemporal patterns of whisker-object contact. *PLoS Computational Biology*, 7(4), 2011.
- [32] Moritz Von Heimendahl, Pavel M Itskov, Ehsan Arabzadeh, and Mathew E Diamond. Neuronal activity in rat barrel cortex underlying texture discrimination. *PLoS Biol*, 5(11):e305, 2007.
- [33] Wikipedia. Bullet (software) — wikipedia, the free encyclopedia, 2016. [Online; accessed 19-October-2016].
- [34] Daniel L K 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 of the United States of America*, 111(23):8619–24, jun 2014.
- [35] Daniel LK Yamins and James J DiCarlo. Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience*, 19(3):356–365, 2016.
- [36] Chunxiu Yu, Dori Derdikman, Sebastian Haidarliu, and Ehud Ahissar. Parallel thalamic pathways for whisking and touch signals in the rat. *PLoS Biol*, 4(5):e124, 2006.

Supplementary information

1 Dataset generation

1.1 ShapeNet category refinement and object preprocessing

Starting from the original 55 categories in ShapeNet [2], we further split them into 117 categories based on the information provided by ShapeNet. For every object in ShapeNet, beside the label from 55 categories, a tree of synsets in WordNet [8] is also provided where the object will be described in more and more detail. For example, one tree could include "artifact", "instrumentation", "transport", "vehicle", "craft", "vessel", and finally "ship". The deeper for one synset in the tree, more refined that synset will be. The label provided is usually not the deepest synset in the tree, which gives rise for refinement. Therefore, we first regrouped the objects in ShapeNet using their deepest synset in synset tree. Then we manually combined some of them as they represent the same thing, such as "liquid crystal display" and "computer monitor". In order to get a dataset with balanced category distribution, we dropped the subcategories containing less than 30 objects, which left 117 subcategories (a full list of the categories can be found in Table 4). Furthermore, with the aim of including roughly 10000 objects with balanced category distribution, we first sorted the categories by the number of objects contained and then sampled objects from smallest category to largest category sequentially. For every sampling, we would first multiply the number of objects in this category with the number of categories left (c_l). If the result was smaller than 10000 minus the number of objects already sampled (n_a), we would just take all of objects in this category. Otherwise, we would randomly sample $(10000 - n_a)/c_l$ objects from each category left. Finally, we got 9981 objects in 117 categories, with most categories containing 91 objects and smallest category containing 41 objects.

The correct collision simulation in Bullet [9] requires the object to be composed of several convex shapes. Therefore, we used V-HACD [7] to decompose every object. The decomposition is done to each object with "resolution" parameter being 500000, "maxNumVerticesPerCH" being 64, and other parameters being default.

1.2 Sweep simulation procedure

Two technical details about sweep simulatin would be described in this section, including rescaling of the object and the adjusting of the position to avoid

collision with fixed whisker bases.

For each object, we first computed a cuboid cover for the whole shape by getting the maximal and minimal values of points on the object in three dimensions. The longest distance in the cuboid cover would then be used as the current scale of this object. And scaling factor was computed by dividing the desired scale by the longest distance.

After rescaling, the object was moved to make the center of cuboid cover to be placed at the desired position and then the nearest point on the object to the center of whisker array was computed. The object was moved again to make the nearest point to be placed at the desired position. To avoid collision, we sampled points on the surfaces of cuboid cover and then computed the trajectories of these points. We would move the object to right to make the nearest distance from every fixed base to every trajectory larger than 4mm.

1.3 Basic control experiment

We generated 24000 independent sweeps for each version of control dataset, meaning 12000 for each category. For each dataset, we split the dataset equally to three parts containing 8000 sweeps overall. We then took 7000 of them as training dataset and the left as testing dataset. To reduce the over-fitting, we randomly sampled 1000 data from input data and the sampling remained the same for all control datasets. The number of units sampled has been searched to make sure that performance could not be added with more units. A LinearSVM is then trained to do the classification with a grid search of parameters. The standard deviation shown on the figure was based on the performances across three splits in one control dataset.

2 Model training and visualization

2.1 Training procedure

We searched the learning parameter space including the choices of learning algorithm, learning rates, and decay rates. Among them, training using Adagrad [3] with a initial learning rates at 0.003 and batch size of 128 gave the best result, which was then applied to the training of all networks. The learning rate would not be changed to 0.0015 until the performance on validation saturated. Ten more epoches would be run after changing the learning rate to 0.0015. The reported performances were the best performances on validation during the whole training procedure.

2.2 Model structures

To make the description of model structures easier, we would use $\text{conv}(\textit{filter size}, \textit{number of filters})$, $\text{pool}(\textit{filter size}, \textit{stride})$, and $\text{fc}(\textit{dimension of outputs})$ to represent convolution layer, pooling layer, and fully connected layer respectively. The size of stride in convolution layers is set to be 1×1 for all networks and the

pooling used in our networks is always max-pooling. For example, $\text{conv}(3 \times 3, 384)$ represents a convolution layer with filter size of 3×3 , stride of 1×1 , and 384 filters.

We used Xavier initialization [5] for the weight matrix in convolution layers while the bias was all initialized to 1. The weight matrix in fully connected layer was initialized using truncated normal distribution with mean being 0 and standard deviation being 0.01, where values deviated from mean for more than two standard deviations would be discarded and redrawn. The bias in fully connected layer was initialized to 0.1. A ReLu layer would always be added after convolution or fully connected layer. And for the fully connected layers, we would use a dropout of rate 0.5 during training. [6]

Simultaneous Spatiotemporal Integration. This family of networks (family S) will usually have several convolution layers followed by fully connected layers. The convolution is applied to both the temporal and spatial dimension.

In Table 1, we showed the structure of all networks in this family, corresponding to the name used in the main paper. The output of previous layer will serve as the input to the next layer. And the layers before fc_combine are shared across three sweeps. For example, in model "S", conv1 to fc7 are shared across sweeps while the outputs of fc7 will be concatenated together as the input to fc_combine8. Besides, S_rand is the same model with S_more except that the weights of conv1 to fc7 there are not trained. Only fc_combine8 is trained to prove that training of weights is necessary.

Separately, we trained another network based on S_more combining 12 sweeps rather than 3 sweeps. The model shared the same structure as S_more. But fc_combine8 there would combine information from 12 sweeps, which means that the input to fc_combine8 is a vector of 1024×12 . The performances of this network surpassed all other networks, with top1 being 0.20 and top5 being 0.53. This means that our network could utilize information from 12 sweeps to help finish the task.

Separate Spatial and Temporal Integration. The network structures for those two families are shown in Table 2. For the "Temporal-Spatial" family, temporal convolution is applied to the inputs first for each whisker using shared weights and then "spatial regroup" in Table 2 means that the outputs from previous layer for each whisker will be grouped according to the spatial position of each whisker in 5×7 grid, with vacant positions filled by zeros. For the "Spatial-Temporal" family, the input is first split into 22 vectors on temporal dimension and each vectore is further reshaped into 2D spatial grid, which means the final shape of each vector would be $5 \times 7 \times 90$. Spatial convolution networks with shared weights will be applied to each vector. After that, the "temporal concatenating" in Table 2 means that the outputs from spatial networks will first be reshaped to one dimension vector and then be concatenated in the temporal dimension to form a new input for further processing. Then temporal convolution will be applied to the new input.

Recurrent Neural Networks with Skip and Feedback Connections. Similar to "Spatial-Temporal" family, the input is first split and reshaped into 22 vectors of size $5 \times 7 \times 90$. And then the vectors are fed into the network one

by one in order of time. The structures of networks in this family are shown in Table 3 with additional edges. The "RNN_lstm" and "RNN_gru" is just adding LSTM/GRU between fc8 and fc_combine9 with number hidden units being 512.

2.3 Visualization related

Calculating category and object level RDMs. After the network training finished, we took the network with best performances in each family and generated the responses of each layer on validation set. For category RDM, we calculated the category mean of responses for each layer and then calculate the RDM based on that using Pearson correlation coefficient between each category. For object level RDM, the average was computed for each object for each layer and then the RDM was calculated from that using Pearson correlation coefficient between each object.

Object grouping. The categories are regrouped using category level RDM to show that the categorization is reasonable. The group is shown in the confusion matrix visualization in main paper, with 117 categories into 10 groups. Specifically, we took the RDM of the highest hidden layer in our best network ("TS") and used affinity propagation to group them [4] with parameter "damping" being 0.9. The results are shown in Table 4.

MDS embedding of RDMs. In order to estimate variance caused by initialization of the DNNs, we trained 5 networks for both "TS" model and "RNN_fdb" model from different initializations. We then generated the responses of these 10 models on validation set and computed object level RDMs based on the responses. For "RNN_fdb", we take the responses of last time point of each layer to compute the RDMs. We also tried RDMs computed from averages and concatenations of responses from last half of time points, the results are similar.

After calculating RDMs, we computed Pearson correlation coefficient between RDMs to get a distance matrix between them. Then we used MDS algorithm for to get a 2D embedding [1].

References

- [1] Ingwer Borg and Patrick JF Groenen. *Modern multidimensional scaling: Theory and applications*. Springer Science & Business Media, 2005.
- [2] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. *ArXiv*, 2015.
- [3] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- [4] Brendan J Frey and Delbert Dueck. Clustering by passing messages between data points. *science*, 315(5814):972–976, 2007.
- [5] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Aistats*, volume 9, pages 249–256, 2010.

- [6] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet Classification with Deep Convolutional Neural Networks. *Advances In Neural Information Processing Systems*, pages 1–9, 2012.
- [7] Khaled Mamou and Faouzi Ghorbel. A simple and efficient approach for 3d mesh approximate convex decomposition. In *Image Processing (ICIP), 2009 16th IEEE International Conference on*, pages 3501–3504. IEEE, 2009.
- [8] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [9] Wikipedia. Bullet (software) — wikipedia, the free encyclopedia, 2016. [Online; accessed 19-October-2016].

Table 1: Network structure for S family

Name	Structure
S	conv1(9×3 , 96), pool1(3×1 , 3×1), conv2(3×3 , 128), pool2(3×3 , 2×2), conv3(3×3 , 192), conv4(3×3 , 192), conv5(3×3 , 128), pool5(3×3 , 2×2), fc6(2048), fc7(1024), fc_combine8(117)
S_more	conv1(9×3 , 96), pool1(3×1 , 3×1), conv2(3×3 , 256), pool2(3×3 , 2×2), conv3(3×3 , 384), conv4(3×3 , 384), conv5(3×3 , 256), pool5(3×3 , 2×2), fc6(4096), fc7(1024), fc_combine8(117)
S_few	conv1(9×3 , 32), pool1(3×1 , 3×1), conv2(3×3 , 64), pool2(3×3 , 2×2), conv3(3×3 , 96), conv4(3×3 , 96), conv5(3×3 , 64), pool5(3×3 , 2×2), fc6(256), fc7(128), fc_combine8(117)
S_4c2f	conv1(9×3 , 96), pool1(3×1 , 3×1), conv2(3×3 , 128), pool2(3×3 , 2×2), conv3(3×3 , 192), conv4(3×3 , 128), pool4(3×3 , 2×2), fc5(2048), fc6(1024), fc_combine7(117)
S_3c2f	conv1(9×3 , 96), pool1(3×1 , 3×1), conv2(3×3 , 388), pool2(3×3 , 2×2), conv3(3×3 , 128), pool3(3×3 , 2×2), fc4(2048), fc5(1024), fc_combine6(117)
S_2c2f	conv1(9×3 , 96), pool1(3×1 , 3×1), conv2(3×3 , 144), pool2(6×6 , 4×4), fc3(2048), fc4(1024), fc_combine5(117)
S_1c2f	conv1(9×3 , 96), pool1(6×3 , 6×2), fc2(896), fc3(512), fc_combine4(117)
S_3c1f	conv1(9×3 , 96), pool1(3×1 , 3×1), conv2(3×3 , 128), pool2(3×3 , 2×2), conv3(3×3 , 192), pool3(3×3 , 2×2), fc4(1532), fc_combine5(117)
S_2c1f	conv1(9×3 , 96), pool1(3×1 , 3×1), conv2(3×3 , 128), pool2(3×3 , 2×2), fc3(708), fc_combine4(117)
S_3c0f	conv1(9×3 , 72), pool1(3×1 , 3×1), conv2(3×3 , 144), pool2(3×3 , 2×2), fc_combine3(117)
S_2c0f	conv1(9×3 , 72), pool1(3×1 , 3×1), fc_combine2(117)
S_1c0f	fc_combine1(117)
S_3D	conv1($9 \times 2 \times 2$, 96), pool1($4 \times 1 \times 1$, $4 \times 1 \times 1$), conv2($3 \times 2 \times 2$, 256), pool2($3 \times 1 \times 1$, $2 \times 1 \times 1$), conv3($3 \times 2 \times 2$, 384), conv4($3 \times 2 \times 2$, 384), conv5($3 \times 2 \times 2$, 256), pool5($3 \times 3 \times 3$, $2 \times 2 \times 2$), fc6(4096), fc7(1024), fc_combine8(117)
S_deep	conv1(5×3 , 64), conv2(5×3 , 64), pool2(3×1 , 3×1), conv3(2×2 , 128), conv4(2×2 , 128), pool4(3×3 , 2×2), conv5(3×3 , 192), conv6(3×3 , 192), conv7(3×3 , 192), conv8(3×3 , 192), conv9(3×3 , 128), pool9(3×3 , 2×2), fc10(2048), fc11(1024), fc12(512), fc_combine13(512), fc_combine14(117)

Table 2: Network structure for "Temporal-Spatial" and "Spatial-Temporal" family

Name	Structure
TS	conv1(9, 64), pool1(3, 3), conv2(3, 256), pool2(3, 2), conv3(3, 384), conv4(3, 384), conv5(3, 256), pool5(3, 2), spatial regroup, conv6(1×1 , 896), conv7(1×1 , 512), conv8(3×3 , 384), conv9(3×3 , 384), conv10(3×3 , 256), fc11(2048), fc12(512), fc_combine13(512), fc_combine14(117)
TS_few	conv1(9, 64), pool1(3, 3), conv2(3, 192), pool2(3, 2), conv3(3, 256), conv4(3, 256), conv5(3, 192), pool5(3, 2), spatial regroup, conv6(1×1 , 512), conv7(1×1 , 384), conv8(3×3 , 256), conv9(3×3 , 256), conv10(3×3 , 192), fc11(1024), fc12(512), fc_combine13(512), fc_combine14(117)
ST	conv1(2×2 , 256), conv2(2×2 , 384), conv3(2×2 , 512), conv4(2×2 , 512), conv5(2×2 , 512), conv6(2×2 , 384), temporal concatenating, conv7(1, 1024), conv8(1, 512), conv9(3, 512), conv10(3, 512), conv11(3, 512), conv12(3, 512), pool12(3, 2), fc13(1024), fc_combine14(512), fc_combine15(117)

Table 3: Network structure for Recurrent Neural Networks family

Name	Structure	Additional Edges
RNN	conv1(2×2 , 96), conv2(2×2 , 256), conv3(2×2 , 384), conv4(2×2 , 384), conv5(2×2 , 384), conv6(2×2 , 256), fc7(2048), fc8(1024), fc_combine9(1024), fc_combine10(512), fc_combine11(117)	
RNN_byp	conv1(2×2 , 96), conv2(2×2 , 128), conv3(2×2 , 256), conv4(2×2 , 384), conv5(2×2 , 384), conv6(2×2 , 256), fc7(1024), fc8(512), fc_combine9(1024), fc_combine10(512), fc_combine11(117)	conv1→conv3, conv1→conv4, conv2→conv4, conv2→conv5, conv3→conv5, conv3→conv6, conv4→conv6, conv2→fc7, conv4→fc7
RNN_fdb	conv1(2×2 , 96), conv2(2×2 , 128), conv3(2×2 , 256), conv4(2×2 , 384), conv5(2×2 , 384), conv6(2×2 , 256), fc7(1024), fc8(512), fc_combine9(1024), fc_combine10(512), fc_combine11(117)	conv3→conv2, conv4→conv2, conv4→conv3, conv5→conv3, conv5→conv4, conv6→conv4, conv6→conv5, conv2→fc7, conv4→fc7

Table 4: Results for category regroup based on RDM

Group	Categories
1	table-tennis table, folding chair, grand piano, lawn chair, rocking chair, windsor chair, swivel chair, park bench, armchair, straight chair, chair
2	drafting table, kitchen table, secretary, pool table, piano, berth, worktable, console table, easy chair, laptop, bench, coffee table, desk, table
3	upright, basket, birdhouse, platform bed, vertical file, dishwasher, convertible, monitor, love seat, printer, microwave, washer, file, stove, bookshelf, dresser, tweeter, bathtub, loudspeaker, cabinet, sofa
4	camera, school bus, bag, pendulum clock, mailbox, planter, bus
5	floor lamp, dagger, delta wing, revolver, propeller plane, carbine, knife, sniper rifle, guitar, rifle, airplane, jet
6	ferry, cabin cruiser, sea boat, cruise ship, yacht, wheeled vehicle, pistol, ship, tender, train, boat
7	limousine, ambulance, stock car, roadster, jeep, beach wagon, wine bottle, cruiser, sports car, convertible, bottle, racer, sport utility, sedan, coupe, car
8	cap, soda can, coffee mug, jar, mug, helmet, bowl, pot, ashcan, vase
9	data input device, remote control, pillow, telephone, screen, clock, liquid crystal display, cellular telephone
10	microphone, earphone, sailboat, motorcycle, table lamp, faucet, lamp