

# NBLAST: Rapid, sensitive comparison of neuronal structure and construction of neuron family databases

Marta Costa<sup>1,2</sup>, Aaron D. Ostrovsky<sup>1</sup>, James D. Manton<sup>1</sup>, Steffen Prohaska<sup>1,3</sup>, Gregory S. X. E. Jefferis<sup>1\*</sup>

<sup>1</sup>Neurobiology Division, MRC Laboratory of Molecular Biology, Cambridge, CB2 0QH, UK

<sup>2</sup>Department of Genetics, University of Cambridge, Cambridge, CB2 3EH, UK

<sup>3</sup>Zuse Institute Berlin (ZIB), 14195 Berlin-Dahlem, Germany

*Please note that this preprint is a public draft not a submitted manuscript. It is released in the expectation that it will be useful to our colleagues “as is” and in order to solicit feedback as we finalise this study for peer review and publication. All presented results are considered reliable, but it is likely that the final manuscript will contain additional analysis; for this reason, the authorship of this draft may well differ from our final submission. There may be some omissions in the Introduction and Discussion sections of the manuscript, for example with respect to citations to the literature (suggestions are welcome for papers we may have missed). Please note that current versions of the open source software described in the manuscript are available by following links in the Experimental Procedures also summarised at <http://jefferislab.org/si/nblast>. Processed data derived from raw data generously made publicly available by third parties (primarily flycircuit.tw) will be made available at least by the time this paper is accepted, hopefully rather sooner; please contact Greg for details. We welcome feedback, queries and suggestions on any aspect of the manuscript, code or data to [jefferis@mrc-lmb.cam.ac.uk](mailto:jefferis@mrc-lmb.cam.ac.uk).*

## Abstract

Efforts to map neural circuits from model organisms including flies and mice are now generating multi-terabyte datasets of 10,000s of labelled neurons. Technologies such as dense EM based reconstruction, and sparse/multicolor labeling with image registration allow neurons to be embedded within the spatial context of a circuit or a whole brain. These ever-expanding data demand new computational tools to search, organize and navigate neurons. We present a simple, but fast and sensitive, algorithm, NBLAST, for measuring pairwise neuronal similarity by position and local geometry. Inspired by the BLAST algorithm for biological sequence data, NBLAST decomposes a query and target neuron into short segments. Each matched segment pair is scored using a log-likelihood ratio scoring matrix empirically defined by the statistics of real matches and non-matches in the data.

We demonstrate the application of a reference implementation of NBLAST to a dataset of 16,129 single *Drosophila* neurons. NBLAST scores are sensitive enough to distinguish 1) two images of the same neuron, 2) two neurons of the same identified neuronal type 3) two neurons of very closely related types. We demonstrate that NBLAST scores can be used to identify neuronal types, such as olfactory projection neurons, with reliability that matches or exceeds expert annotation in a fraction of the time. We also show that clustering using appropriately normalized NBLAST scores can reveal classic morphological types as well as identify unpublished classes. We carry out detailed analysis of a number of neuronal classes including Kenyon cells, olfactory

and visual projection neurons, auditory, and male-specific P1 neurons. This identifies many new neuronal types and reveals unreported features of topographic organization. Finally we provide a complete clustering of the 16,129 neurons in the test dataset into 1,052 clusters of highly related neurons. These clusters are then organized into superclusters, enabling both exploration of the dataset and the matching of individual clusters to morphological types in the literature. Finally NBLAST queries can be used to identify candidate neurons matching neurite tracts with transgene expression pattern.

We provide a general purpose open source toolbox that implements construction of score matrices, the core pairwise scoring algorithm, *de novo* and precomputed database search and clustering along with complete source code and data for the analyses in this paper. The neuronal families can also be queried online through [virtualflybrain.org](http://virtualflybrain.org) and visualized in interactive 3D at [jefferislab.org](http://jefferislab.org).

## Introduction

Understanding brain function at the level of circuits of interconnected single neurons is one of the major challenges of neuroscience. One of the basic activities in the study of neuronal circuits is to correlate the functional properties and behavioral relevance of neurons with their cell type. There is no universally accepted definition of neuronal cell type, but features including morphology, position within the nervous system, genetic markers and connectivity are typically key descriptors (Migliore and Shepherd, 2005; Bota and Swanson, 2007; Rowe and Stone, 1976). Despite this ambiguity, cell type remains a key and almost universal abstraction for collecting and comparing experimental results from different groups. Since the morphology and position of neurons strongly constrain (and indeed are partially defined by) connectivity, they have been mainstays of studies of circuit organization for over a century. Classic experimental approaches to neuronal morphology include the Golgi method made famous by Cajal, microinjection, and filling of cells during whole cell recording. More recently, these have been supplemented by genetic approaches to sparse and combinatorial neuronal labeling that have enabled increasingly large-scale characterization of single neuron morphology (Jefferis and Livet, 2012). These approaches have been applied in a focussed way to characterize the neurons associated with specific circuits (Marin et al., 2002; Wong et al., 2002; Badea et al., 2003; Jefferis et al., 2007; Lin et al., 2007; Sömböl et al., 2014; Miyasaka et al., 2014).

Classic approaches to map the position of neurons within the brain have depended on visual comparison with anatomical features, often revealed by a general counterstain; these approaches have been most successful when applied to laminar organization as typified by the mammalian retina (Badea and Nathans, 2004; Coombs et al., 2006; Kong et al., 2005; Sömböl et al., 2014), fly optic lobe (Fischbach and Dittrich, 1989; Morante and Desplan, 2008) or cerebellum (Cajal and Azoulay y, 1911). Recently, the use of 3D light microscopy and image registration has enabled direct fusion of datasets to generate digital 3D atlases that can be used to generate very specific testable hypotheses about circuit connectivity, organization and function at large scales (Sunkin et al., 2013; Oh et al., 2014; Zingg et al., 2014; Rybak et al., 2010; El Jundi et al., 2009). The largest study to date using these approaches is that of Chiang et al. (2011), who combined genetic mosaic approaches and image registration to produce an atlas of over 16,000 single cell morphologies embedded within a standard *Drosophila* brain at <http://flycircuit.tw>.

There is an increasing appreciation that cataloguing cell types in an accessible and quantitative fashion can be a driver for circuit research, helping to reveal organizational principles and facilitating the integration of morphological and functional data from different groups. For example, attempts to categorize cortical interneurons by morphology, physiology and gene expression patterns in the last decade have substantially assisted ongoing research into their circuit function (Petilla Interneuron Nomenclature Group et al., 2008; Nelson et al., 2006; Kepecs and Fishell, 2014). Indeed, obtaining a catalogue of neurons in the brain is the first explicit goal of the recently announced US BRAIN Project (<http://www.nih.gov/science/brain/11252013-Interim-Report-Final.pdf>).

The existence of large databases of neuronal morphologies now exceeding 10,000 neurons (Parekh and Ascoli, 2013) presents an acute practical challenge: to find and organize related neurons. However, this challenge also represents an opportunity: quantitative classification of neuronal morphology, especially within the context of template brains, could provide a very substantial contribution to the problem of cell type. A key requirement is a tool to easily and sensitively compute the morphological similarity of neurons within and between datasets. This problem has clear analogies to the issue of

comparing sequence and protein data. The explosion of biological sequence information in the late 80s and early 90s motivated the development and refinement of sequence similarity tools such as FASTA (Pearson and Lipman, 1988) and BLAST (Altschul et al., 1990). These algorithms could be used for pairwise similarity scoring and alignment, and for rapid querying of sequence databases. Similar approaches also provided the basis for generating non-redundant protein databases and the construction of hierarchically organized databases of protein families (Sonnhammer et al., 1997; Joseph et al., 2014)

Neuronal morphology can be represented in digital form as a collection of one or more rooted trees i.e. a directed graph structure embedded in 3D space. Unfortunately, this space has not typically been a defined atlas space, but simply the physical space of the microscope system used to reconstruct the neuron. For this reason, databases such as [NeuroMorpho.org](#) (Ascoli et al., 2007) do not include precise information about the position of neurons within a template brain.

The comparison of arbitrary graph structures is a well-developed problem area (Conte et al., 2004) and there has been some progress in applying these ideas to neurons. For example Mayerich et al. (2012) developed an algorithm, NetMets, to compare neuronal reconstructions with a ground truth. However, since this algorithm requires a (manually) defined ground truth, it is not a general purpose neuronal similarity algorithm suitable for application to large databases. Basu et al. (2011) have proposed the path2path algorithm, a geometric measure of similarity that decomposes neuronal trees into a family of unbranched paths. Finally, Cardona et al. (2010) have decomposed unbranched neurites into sequences of vectors and used a dynamic programming algorithm to generate an optimal alignment between arbitrary 3D curves. They applied this approach to the secondary neuroblast lineages in *Drosophila*, achieving very high classification accuracy against a database of a few hundred traced structures. Although this algorithm could be modified for use with single neurons, it treats each unbranched neuronal segment separately and there is not a natural way to handle trees with many segments.

Tracing even single neurons from image data still remains a challenging problem that has not been fully automated (Brown et al., 2011), while more complex expression patterns cannot be directly represented in a simple fully-connected binary tree as typically used to represent neurons. There are therefore advantages to strategies that do not attempt to represent neuronal structures in a connected graph structure but instead emphasize local geometry and position. For example, we have previously developed an image segmentation pipeline that represents complex neuronal expression patterns (consisting of up to 100 neurons) as point clouds along with tangent vectors defining the local heading of the neuron. This resulted in a supervised learning approach to the challenging problem of recognizing groups of lineage-related neurons (Masse et al., 2012). Ganglberger et al. (2014) have recently applied a somewhat related approach by using local intensity gradients to define vector fields for expression pattern images in which the vectors track the local structure of neurons. Two images can then be compared by calculating the inner product of the two vector fields to produce a distance measure. While this method can be applied to expression pattern data that have not been segmented, the lack of precision inherent in a slowly varying vector field makes it comparatively insensitive when applied to single neurons, and the processed data are very large – of the same order as the original image i.e. 100 MB – resulting in low throughput.

Combining the data representation that we developed in Masse et al. (2012) with a very large single neuron dataset (Chiang et al., 2011) allowed us to test and validate a new algorithm that is both extremely sensitive and very fast. The search data structure is compact (<2 GB for 16,000 neurons) and typical pairwise search times are < 5 ms on a single processor. The algorithm's scoring parameters are



defined by the statistics of matches in the data rather than by expert intuition. We find this to be a very effective tool for neuron search and unsupervised clustering of neurons. We first describe the detailed motivation and implementation of this algorithm and validate its use in database search. We then show that it can be used to identify classic neuronal types in *Drosophila* with a sensitivity matching that of domain experts, in a fraction of the time of manual classification. We go on to show that this procedure can identify new neuronal types and reveal undescribed features of topographic organization. Finally, we apply our method to 16,129 neurons from the FlyCircuit dataset, reducing this to a non-redundant set of 1,052 morphological clusters. Manual evaluation of a subset of these clusters show that they closely match expert definition of cell types. These clusters therefore represent a preliminary global cell type classification for the *Drosophila* brain. They are organized into a supercluster hierarchy that enables easy exploration at different levels, simplifying exploration of the full diversity of this dataset and comparison of related cell types and families.

## Results

### Algorithm

Our principal design goals were to develop a neuron similarity algorithm that included aspects of both spatial location (within a brain or brain region) and neuronal branching pattern, and that was both extremely sensitive and very fast. The applications that we had in mind were the searching of large databases of neurons (10,000-100,000 neurons), clustering of neurons into families by calculating all-against-all similarity matrices, and the efficient organization and navigation of datasets of this size. Our first application would be data acquired in *Drosophila*, where previous studies using image registration have shown a high degree of spatial stereotypy (standard deviation of landmarks  $\sim 2.5$   $\mu\text{m}$  in each axis for a brain of 600  $\mu\text{m}$  in its longest axis, [Jefferis et al., 2007](#)). We therefore expected to use data that had been co-registered rather than calculating similarity using features of the neuron that are independent of absolute spatial location.

One of the first design decisions that we had to make was whether to develop a direct pairwise comparison algorithm or to use a form of dimension reduction to map neuronal structure into a lower-dimensional space. The major advantage of the latter approach is that the similarity between neurons can be computed directly and almost instantaneously in the low dimensional space. However, the construction of a suitable embedding function either requires existing knowledge of neuronal similarity (likely supplied by experts in the form of large amounts of training data), huge amounts of unlabeled data that enable direct learning of features (e.g. [Le et al., 2012](#)), or a strategy based on successful extraction of key image features.

A number of considerations made us favor the approach of direct pairwise comparison. First, we suspected that it would be possible to make a more sensitive algorithm by working with the original data. Second, the amounts of image data available did not seem large enough to avoid a requirement for extensive labeled training data. Third, we reasoned that our own intuitions about neuronal similarity could be better expressed in the original physical space of the neuron than in a low dimensional embedding. Our own exploratory analysis (SP, GSXEJ unpublished observations) confirmed that constructing a sensitive metric of this sort is challenging.

The selection of a pairwise similarity metric meant that we had to give particularly careful consideration to performance issues in the design phase. We set two practical performance targets: 1) being

able to carry out searches of a single neuron against a database of 10,000 neurons in less than a minute on a simple desktop or laptop computer. 2) Being able to complete all-against-all searches for 10,000 neurons ( $10^8$  comparisons) in  $< 1$  day on a powerful desktop computer. These targets meant that each elementary comparison operation should take around 5 ms or less. Image pre-processing carried out once per neuron would therefore be a good investment if it reduced the time taken for each pairwise comparison. These considerations prompted us to generate a spatially registered, compact representation of each neuron as a separate pre-processing step for each neuron, rather than develop an algorithm that simultaneously solved both the spatial alignment and similarity problem.

The starting point for our algorithm is a representation in which neuron structures have been reduced to segments, represented as a location and an associated local tangent vector. This retains some local geometry but does not attempt to capture the topology of the neuron's branching structure. We have found that a simplified representation of this sort can be constructed for image data that would not permit automated reconstructions. In order to prepare data of this sort in quantity, we developed an image processing pipeline summarized in Figure 1A and detailed in Experimental Procedures. Briefly, brain images from the FlyCircuit dataset (Chiang et al., 2011) were subjected to non-rigid image registration (Jefferis et al., 2007) against a newly constructed intersex template brain. Neuron images were thresholded and subjected to a 3D skeletonization procedure (Lee et al., 1994) implemented in Fiji (Schindelin et al., 2012). These thresholded images were then converted to the point and tangent vector representation (Masse et al., 2012) using our R package `nat` (Jefferis and Manton, 2014); the tangent vector (i.e. the local heading) of the neuron at each point was computed as the first eigenvector of a singular value decomposition (SVD) of the point and its 5 nearest neighbors.

After this pre-processing, data could then be loaded into R and plotted in 3D and analyzed using `nat` (Figure 1B). Neurons were represented by a median of 1,070 points/vectors each, occupying 61.5 kB each or 1.8 GB for all 16,129 neurons, therefore fitting comfortably into main memory on a standard personal computer. Since the fly brain is almost completely symmetric, but neurons were labelled randomly in both hemispheres, we opted to map all neurons to the left hemisphere (defined primarily by cell body location, see Experimental Procedures and Figure 1B) using a non-rigid mirroring procedure (Manton et al., 2014). In this way identified cell types labeled on opposite sides of the brain should be co-aligned.

With a database of aligned neurons in an appropriate representation, we were then able to calculate NBLAST pairwise similarity scores. One neuron is designated the query and the other the target. For each query segment (defined by a midpoint and tangent vector) the nearest neighbor (using straightforward Euclidean distance) is identified in the target neuron (Figure 1C–D). A score for the segment pair is calculated as a function of two measurements:  $d$ , the distance between the matched segments (indexed by  $i$ ) and  $|\vec{u}_i \cdot \vec{v}_i|$ , the absolute dot product of the two tangent vectors; the absolute dot product is used because the orientation of the tangent vectors has no meaning in our data representation (Figure 1C). The scores are then summed over each segment pair to give a raw score,  $S$ :

$$S(\text{query}, \text{target}) = \sum_{i=1}^n f(d_i, |\vec{u}_i \cdot \vec{v}_i|). \quad (1)$$

The question then becomes: what is an appropriate function  $f(d_i, |\vec{u}_i \cdot \vec{v}_i|)$ ? We initially experimented with a function based on expert intuition:

$$f = \sqrt{|\vec{u}_i \cdot \vec{v}_i| \exp(-d_i^2/2\sigma^2)}. \quad (2)$$

This includes a negative exponential function of distance (related to the normal distribution), with a free parameter  $\sigma$  based on our previous estimates of the variability in position within the fly brain of landmarks after registration (Jefferis et al., 2007; Yu et al., 2010b) set to 3  $\mu\text{m}$ . Although this provided a useful starting point, we were unhappy with a scoring system that required parameters to be specified rather than derived empirically from data. We therefore developed an alternative approach inspired by the scoring system of the BLAST algorithm (Altschul et al., 1990). For each segment pair we defined the score as the log probability ratio:

$$f = \log_2 \frac{p_{\text{match}}}{p_{\text{rand}}}, \quad (3)$$

i.e. the probability that the segment pair was derived from a pair of neurons of the same type, versus a pair of unrelated neurons. We could then define  $p_{\text{match}}$  empirically by finding the joint distribution of  $d$  and  $|\vec{u}_i \cdot \vec{v}_i|$  for pairs of neurons of the same type (Figure 1E–G). For our default scoring matrix, we used a set of 150 olfactory projection neurons innervating the same glomerulus, which can be unambiguously defined as the same neuronal type (Figure 1F).  $p_{\text{rand}}$  was calculated simply by drawing 5,000 random pairs of neurons from the database, assuming that the large majority of such pairs are not morphologically related neurons. The joint distributions were calculated using 10 bins for the absolute dot product and 21 bins for the distance to give two 21 row  $\times$  10 column matrices. The 2D histograms were then normalized to convert them to probabilities and the log ratio then defined the final scoring matrix (Figure 1G). Plotting the scoring matrix emphasizes the strong distance dependence of the score but also shows that for segment pairs closer than  $\sim 10 \mu\text{m}$ , the logarithm of the odds score increases markedly as the absolute dot product moves from 0 to 1 (Figure 1H).

We implemented the NBLAST algorithm as an R package ([nat.nblast](#)) that immediately enables pairwise queries, searches of a single query neuron against a database of target neurons (Figure 2A) and all-by-all searches. The median number of segments for neurons in the test dataset was 1,070. Run times on a single core of a laptop computer for neurons of this size were 2 ms per comparison or 30 s for all 16,129 neurons. In order to enable clustering of neurons on the fly, we also pre-computed an all-by-all similarity matrix for all 16,129 neurons ( $2.6 \times 10^8$  scores). This occupied 1.0 GB on disk, uncompressed. If desired, subregions of the matrix could be rapidly loaded into memory on demand rather than retaining the entire matrix in memory.

## Searching with NBLAST

A good similarity algorithm should be sensitive enough to reveal identical neurons with certainty, while having the specificity to ensure that all high scoring results are relevant hits. We used the full FlyCircuit dataset to validate the performance of NBLAST. We started by carrying out qualitative assessment of searches from single query neurons. Manual inspection of the ranked results by expert anatomists indicated that NBLAST scores are a very good indicator of similarity: highly scoring neurons were rated as very similar and lower scoring neurons as less similar to the query.

As an example we present a search using an auditory interneuron, [fru-M-300198](#), as query (Figure 2B–D). Ranking query results by their NBLAST score, the first listed object is the query neuron

itself (since it is present in the database). The next neuron (**fru-M-300174**), the top hit for the search, overlaps with the query (Figure 2B’); indeed the match is so close that we suspected that they were the same specimen. Further investigation revealed that these “identical twins”, both derived from the same raw confocal image. The next 8 hits are also very similar to the query but are clearly distinct specimens, having small differences in position, length and neurite branching that are typical of sister neurons of the same type (Figure 2B’’).

When we examined the NBLAST scores, we found that the top hit was clearly an outlier; its NBLAST score was 9,471, or 96.1 % of the self-match score of the query against itself (which is the maximum possible score). The small difference in score is likely due to subtle differences in the segmentations of the two neurons. The second hit, although still very similar to the query neuron, has a considerably lower score. A histogram of all scores reveals that only a minority of hits (3 %) have a score above 0 (Figure 2C–D), with the top 8 hits having a score over 7,000. This score of 0 represents a natural cutoff for the NBLAST algorithm since it means that, on average, segments from this neuron have a similarity level that is equally likely to have arisen from a random pair of neurons in the database as a pair of neurons of the same type.

Dividing the histogram of scores into bins and plotting these corresponding groups of neurons shows that there is a clear gradation of similarity, with lower scores corresponding to neurons that are less alike (Figure 2D’). In addition, it shows that only the highest scoring neurons (groups I, II and part of III) are what we would classify as neurons of the same class. These results show that the NBLAST algorithm is very sensitive, with very small differences in neuron morphology between two neurons being reflected in their score.

Although raw NBLAST scores correctly distinguishes between similar and non-similar neurons and can highlight strong partial matches, they are not comparable from one query neuron to the next, since the score depends on neuron size and number of segments. This is also true when reversing the identity of query and target neurons and can confound results when searches are done between neurons of very different sizes (Figure 2E). For example, a search with a large neuron as query and a smaller one as target (pair 1) will have a very low forward score, because the large neuron has many segments that are unmatched, but a high reverse score, since most of target will match part of the query. This reverse score might in fact be higher than a search between two similar neurons (pair 2), erroneously suggesting that the neurons in pair 1 are more similar than in pair 2. One approach to correct for this is to calculate normalized scores by dividing the raw score for each pair by the self-score for the query. Although normalized scores are comparable between pairs, the issue of very unequal forward and reverse scores between large and small neurons remains. One simple strategy to address this issue is to calculate both forward and reverse scores then take the mean of two scores (mean score). In our example (Figure 2E) the mean score for neuron pair 1 is lower than for pair 2, showing that mean scores accurately reflect the similarity between pairs of neurons. We performed a similar analysis to Figure 2D–D’ using mean scores (Figure S2); this eliminated some matches between the query neuron and large neurons with extensive but relatively non-specific arbors that were not considered by biologists to be relevant matches.

During analysis of the FlyCircuit dataset, we sporadically found other cases of apparent human error, where two neurons in the database turned out to derive from the same physical neuron (Figure S1). If the NBLAST algorithm is sensitive enough, it should be able to find other instances of this, for example, duplicated raw images. In order to test this, we collected the top ranking pairs from the all-by-all matrix, i.e. the top hit for each neuron. We analyzed the distribution of forward (Fig-

ure 2F) and reverse scores (data not shown) for all top hits. The distribution is skewed with a long tail of low scores, a mean of 0.58, a peak (i.e. mode) of 0.735, and then falls precipitously with a small tail amounting to 1 % of all top hits, that have anomalously high scores  $> 0.8$ . Given this distribution, we further examined neuron pairs with both forward and reverse scores  $> 0.8$ . Using the image meta-data and looking at the original stacks, we classified 72 pairs of neurons into 4 different groups. From highest to lowest predicted similarity the groups are: same segmentation, when a segmented image of a neuron has been duplicated (Figure S1A); same raw image, corresponding to a different segmentation of the same neuron (Figure 2B'); same specimen, when two images are from the same brain but not from the same confocal image (Figure S1B); and different specimen, when two neurons are from different brains, and the high score implies a very high level of similarity, suggesting that these neurons are likely to be of the same type. A plot of the scores for these pairs shows that the distribution of scores for the different categories follows the predicted similarity (Figure 2G). The highest scores are for the same segmentation group, whereas the lowest corresponds to different specimen pairs. The scores for the same raw image and same specimen groups are mixed, although the scores of the same raw image pairs are slightly higher than the same specimen group pairs. These results underline the high sensitivity of the NBLAST algorithm to small differences between neurons.

These results validate the NBLAST algorithm as an accurate tool for searching for similar neurons. Our method of image preprocessing is able to capture and retain important morphological features of a neuron. The NBLAST algorithm is very sensitive to small differences in morphology, with the highest scores in a search correctly identifying the most similar neurons to the query. In addition, the mean score provides a directly comparable measure of similarity between pairs of neurons.

## Clustering NBLAST scores can identify Kenyon cell classes

We have shown that NBLAST can quantify the similarity of a query neuron to a database of target neurons. We now investigated whether NBLAST scores can be used to cluster neurons by structure and position, potentially revealing functional classes. This is directly analogous to clustering proteins by sequence similarity scores, an approach that can be an invaluable starting point for functional investigation (e.g. Manning et al., 2002). We decided to begin our investigation of this issue with Kenyon cells (KCs). These are the intrinsic neurons of the mushroom body neuropil and are probably the most extensively studied category of neuron in the fly brain given their key role in memory formation and retrieval (reviewed in Kahsai and Zars, 2011).

There are around 2,500 KCs in each mushroom body, and they form the medial lobe, consisting of the  $\gamma$ ,  $\beta'$  and  $\beta$  lobes, the vertical lobe, consisting of the  $\alpha$  and  $\alpha'$  lobes, the calyx, where dendrites are and around which cell bodies are positioned, and the peduncle, formed by the anterior projection of the axons before joining the lobes. The main synaptic input region is the calyx, whereas output connections are located in the lobes. Three main classes of KCs and a few subclasses of neurons, which differ in morphology and birth time are recognized: the  $\gamma$  neurons are the first to be born and innervate only the  $\gamma$  lobe; the  $\alpha'/\beta'$  neurons are generated next and bifurcate at the anterior end of the peduncle and project to the  $\alpha'$  and  $\beta'$  lobes; the  $\alpha/\beta$  neurons are the last to be born, and project to both the  $\alpha$  and  $\beta$  lobes. Four neuroblast clones which differ in their position in the calyx generate the KCs, with each one generating the whole repertoire of neuron types (Lee et al., 1999).

We wanted to investigate if clustering the KCs present in the FlyCircuit dataset based on NBLAST scores would generate groups that reflect the classification of KCs into  $\gamma$ ,  $\alpha/\beta$  and  $\alpha'/\beta'$  neurons. The



first step was to identify and collect all the KCs. In order to do this we performed a forward and reverse search against all neurons using one identified KC ([fru-M-500225](#)), and selected neurons that had both raw scores above  $-2,500$  (2,088 neurons). We performed affinity propagation clustering ([Frey and Dueck, 2007](#)) of these neurons, obtaining 59 clusters, and manually verified each one, resulting in 1,562 neurons being identified as KCs. An additional search for high scorers against these KC exemplars uncovered an extra 102 neurons, bringing the total number of KCs used in our analysis to 1,664, representing 10.3 % of the FlyCircuit dataset.

We performed hierarchical clustering on the KCs, dividing the dendrogram into two groups (Figure 3A). Contrary to expectations, one group contained both the  $\gamma$  and  $\alpha'/\beta'$  neurons, whereas the other group consisted exclusively of  $\alpha/\beta$  neurons (Figure 3B–D), the largest subset in our sample. We separated  $\alpha'/\beta'$  from the  $\gamma$  neurons in a subsequent hierarchical clustering of this group. We performed additional analysis for each of the neuron types.

Hierarchical clustering of the 470  $\gamma$  neurons resulted in a dendrogram which we divided into three groups (I–III) (Figure 3B'). The number of clusters was chosen by visual inspection in order to reveal differences in morphology and organization between the groups. Groups I and III corresponded to the classical  $\gamma$  neurons while group II matched atypical  $\gamma$  neurons. There were differences in neurite positioning in the calyx, from medial to lateral, with group I being the most medial, followed by groups II and III. There were also differences in the gamma lobe, with group II occupying the anterodorsal region, while groups II and III were mostly mixed in the rest of the lobe. A subsequent clustering analysis of the classical  $\gamma$  neurons divided into 4 groups revealed that there were differences between the groups in their medial to lateral position in the calyx. These differences correlated to a certain degree with differences in the dorsal/ventral position of the projections in the  $\gamma$  lobe, with group I, the most medial, being also the most dorsal (Figure 3B''). These observations suggest that the relative position of the projections of classical  $\gamma$  neurons is maintained at the input (calyx) and output sites ( $\gamma$  lobe). We experimented with clustering the classic  $\gamma$  neurons based only on the scores of the segments in the peduncle. The overall organization almost fully recapitulated the positioning of the neurites in the whole neuron analysis (for more information see Figure S3 and ). Thus, the stereotypical organization of the classical  $\gamma$  neurons is maintained throughout the neuropil, from the input to the output sites.

Group II of the  $\gamma$  neurons matched atypical  $\gamma$  neurons, with neurons that extended neurites posteriolaterally in the calyx and projected to the most dorsal region of the  $\gamma$  lobe (Figure 3B'''). Hierarchical clustering of these neurons, resulted in a dendrogram that we divided into 3 groups (a–c). This number of groups isolated the atypical  $\gamma$ d neurons into one group (group a) ([Aso et al., 2009](#)). These neurons extend neurites ventrolaterally at the level of the calyx. The other 2 groups (b, c) correspond to uncharacterized types. Although they project to a similar region in the  $\gamma$  lobe, their dendrites do not extend laterally and their calyx neurites are longer than  $\gamma$ d neurons.

Hierarchical clustering analysis of the  $\alpha'/\beta'$  neurons and separation into 4 groups highlighted the characterized subtypes of  $\alpha'/\beta'$  neurons (Figure 3C–C'). They differ in their anterior/posterior position in the peduncle and  $\beta'$  lobe with three types described -  $\alpha'/\beta'$  anterior,  $\alpha'/\beta'$  medial and  $\alpha'/\beta'$  posterior ([Tanaka et al., 2008](#)). Although we were unable to unambiguously assign a  $\alpha'/\beta'$  subtype to each group i–iv there were clear trends. Neurites of neurons in groups i and iv were more anterior than the other 2 groups (ii, iii) in both the peduncle and  $\beta'$  lobe. These relative positions were not maintained in the calyx, with the two the anterior groups (i, iv) occupying either a medial or a lateral position.

The largest subset of KCs corresponds to  $\alpha/\beta$  neurons (Figure 3D). During the analysis of this group we found 18 neurons that did not correspond to  $\alpha/\beta$  cells, since they innervated either only the  $\beta$  or the  $\alpha$



lobes, and they were removed from the analysis. We performed hierarchical clustering on the remaining 1,091 cells and divided the resulting dendrogram into four groups (1–4) (Figure 3D'), which matched the four neuroblast lineages from which they originate (Zhu et al., 2003). The relative position of the neurites of the four groups within the calyx is somewhat maintained in the peduncle, with the lateral neuroblast clones (group 1 AL; group 2 PL) extending along the dorsolateral peduncle, while the medial clones (group 3 AM; group 4 PM) occupy a more ventromedial region. Hierarchical clustering of each group revealed an expected common organization for all neuroblast clones (Figure 3D'', Figure S3A). For all groups there was a clear distinction between the late born core ( $\alpha/\beta$  core,  $\alpha/\beta$ -c) and early born peripheral neurons ( $\alpha/\beta$  surface,  $\alpha/\beta$ -s). Core neurons are on the inside stratum of the  $\alpha$  lobe. They are also reported to occupy the inside stratum of the peduncle and  $\beta$  lobe (Tanaka et al., 2008). We did not observe this, although the projections of  $\alpha/\beta$  core neurons were ventroposterior to  $\alpha/\beta$  surface neurons in both the peduncle and  $\beta$  lobe. There was also a trend for  $\alpha/\beta$  surface neurons to occupy a more medial position in the calyx in comparison to  $\alpha/\beta$  core ones. A subgroup of group 2 corresponded to the  $\alpha/\beta$  posterior or pioneers neurons ( $\alpha/\beta$ p). The  $\alpha/\beta$ p neurons are the earliest born  $\alpha/\beta$  and they innervate the accessory calyx, run along the surface of the posterior peduncle into the  $\beta$  lobe but stop before reaching the medial tip. A new clustering based on peduncle position of the neuron segments did not recapitulate the relative positions of the calyx neurites for each of the neuroblast clones observed in the whole neuron analysis suggesting that the relative position of the  $\alpha/\beta$  neurons in the peduncle does not completely reflect their stereotypical organization in the calyx (for more information see Figure S3 and ).

In summary, the hierarchical clustering of KCs using the NBLAST scores resolved the neurons into the previously described KC types and some of the subtypes, and isolated uncharacterized subtypes in an extensively studied cell population. In addition, it revealed organizational principles that have been previously described (Tanaka et al., 2008). These observations support our claim that the NBLAST scores retain enough morphological information to accurately search for similar neurons and organize large datasets of related cells.

## Clustering can identify classic cell types: olfactory projection neuron clustering

We have shown that clustering based on NBLAST scores can identify the major classes of Kenyon cells along with a number of sub-classes, some of which have only recently been described, while others may be novel. However the very large number of Kenyon cells along with the variability in their inputs (Caron et al., 2013) means that it is rather unclear what corresponds to an identified cell type, which we take to be the finest classification of neuron in the brain. We therefore analyzed a different neuron family, the olfactory projection neurons that provide input to Kenyon cells, and represent one of the best defined and most intensively studied examples of defined cell types in the fly brain.

Olfactory projection neurons (PNs) transmit information between antennal lobe glomeruli, which receive sensory input, and higher olfactory brain centers, including the mushroom body and the lateral horn (Masse et al., 2009). PNs can be classified on the basis of the number of glomeruli they innervate: poly- and oligoglomerular PNs innervate multiple glomeruli. Uniglomerular PNs (uPNs), whose dendrites innervate just one glomerulus, have been intensively studied and are highly stereotyped in both morphology and developmental origin. They are classified into individual types based on the glomerulus they innervate and the axon tract they follow; these features show fixed relationships with their axonal branching patterns in higher brain centers and their parental neuroblast (Marin et al.,

2002; Jefferis et al., 2001; Wong et al., 2002; Jefferis et al., 2007; Yu et al., 2010a; Tanaka et al., 2012). There are 56 glomeruli in the antennal lobe (Tanaka et al., 2012) and a slightly larger number of uPN types, since some glomeruli have more than one type of post-synaptic projection neuron (Marin et al., 2002; Wong et al., 2002). In some cases (such as embryonic-born PN types of the anterodorsal lineage) it is known that there is just one neuron per hemisphere of a particular PN type (Yu et al., 2010a), while there are 6-7 PN types innervating DA1, one of the largest glomeruli (Jefferis et al., 2004).

We investigated whether the NBLAST algorithm was sensitive enough to capture the very fine level of detail of individual uPN types. We started by manually classifying the 400 uPNs in the Fly-Circuit dataset by glomerulus, neuroblast lineage, and axon tract, using the original image stacks. The definition of the manual gold standard annotations was an iterative process that took several days. The first round accuracy was about 95 %. Numerous discrepancies were revealed by subsequent NBLAST analysis and difficult cases were resolved by discussion between two expert annotators before finalizing assignments. We excluded 3 neurons for which no conclusion could be reached. We found a very large number of DL2 uPNs, 145 DL2d and 37 DL2v, in a total of 397 neurons. Nevertheless, our final set of uPNs broadly represents the total variability of described classes and contains neurons innervating 35 out of 56 different glomeruli (Tanaka et al., 2012), examples of the three main lineage clones (adPNs, lPNs and vPNs) in addition to one bilateral uPN, and neurons that follow each of the three main tracts (medial, mediolateral and lateral antennal lobe tracts).

We then used NBLAST to analyze these neurons, dropping 3 neurons for which registrations failed. We computed NBLAST scores for each uPN against the remaining 16,128 neurons in the test dataset and checked whether the top hit was exactly the same type of PN, a similar uPN or a match to another class of neuron. We restricted our analysis to types with at least two examples in the dataset and to unique pairs (i.e. if PN A was the top hit for PN B and vice versa, we only counted them once) (n=327). We then summarized the number of matches in each category according to their forward and reverse normalized scores (Figure 4A). 97.6 % of top hits were to exactly the same PN class (n=319). Of the remaining 8 neurons, 4 had matches to a uPN innervating a neighboring glomerulus with identical axon projections (DL2d vs DL2v, VM5d vs VM5v) that are challenging even for experts to distinguish (for example of VM5v and VM5d compare groups 2 and 3 in Figure 4C). There were a further four matches to neurons that were not uPNs, in which the top hit was an oligoPN that innervated the same glomerulus as the query. These results show that NBLAST is sensitive enough to capture the level of detail required to distinguish uPN types.

In addition to examining the top NBLAST hits, we also compared how the top 3 hits matched the query type (Figure 4B). We removed DL2 uPNs from the set of neurons analyzed, given their high prevalence in the dataset. For the remaining uPN types with more than three examples (n=190), we collected the top three NBLAST hits. In 99.5 % of cases (i.e. all except one) at least one of top hits matched the query type, decreasing to 96.3 % and 92 % for at least two or all three matches, respectively. For 97.9 % of instances the top hit was the same type as the query.

Given the very high prediction accuracy of neuron type by NBLAST searches, we wondered if an unsupervised clustering based on NBLAST scores would also group them by type. To test this, we clustered uPNs (without DL2 neurons, n=214) and divided the dendrogram at a height of 0.75, as this level was found to be the one at which most groups corresponded to single and unique neuron types (Figure 4C). For types with more than one representative neuron, all neurons co-clustered, with three exceptions: the isolated VM5v neuron in group 8, and the DL1 and DA1 neurons that are split between two groups (groups 12 and 13, and groups 15 and 16, respectively). The cluster organization

also reflects higher level features such as the axon tract / neuroblast of origin, as well as the detailed dendritic position and axon terminal arbor morphology. Thus, unsupervised clustering of uPNs based on NBLAST scores gives an almost perfect neuronal classification.

In conclusion, these results demonstrate that morphological comparison by NBLAST is powerful enough to resolve differences at the finest level of neuronal classification. Furthermore, they suggest that unsupervised clustering by NBLAST scores could help to reveal new neuronal types.

## Clustering can be used to define new cell types

### Visual projection neurons

Visual projection neurons (VPNs) relay information between the optic lobes and the central brain. They are a morphologically diverse group that innervate distinct optic lobe and central brain neuropils, with 44 types already described (Otsuna and Ito, 2006). We explored whether clustering of these neurons based on NBLAST scores would find previously reported neuron classes and identify new ones.

We started by selecting the VPns from all the neurons in the test dataset. We started with the 1,052 exemplars found by affinity propagation clustering of NBLAST scores (Figure 8). We then clustered those exemplars using hierarchical clustering and found that extrinsic and intrinsic optic lobe neurons together formed a distinct “optic lobe” group within this (Figure 8C). We then collected all neurons associated with those “optic lobe” exemplars and calculated the overlap of neurons with each of the standard neuropils defined by Ito et al. (2014) (see Experimental Procedures and Manton et al., 2014 for technical details). This then allowed us to separate neurons by innervation pattern into 3 groups: 1) ipsilateral optic lobe neuropils only (see Figure S6) ipsilateral and central brain neuropils (unilateral VPns, uVPns) or 3) both optic lobes and central brain neuropils (bilateral VPns). This selection procedure resulted in a set of 1,793 uVPns and 72 bilateral VPns.

Hierarchical clustering using NBLAST scores for the uVPns resulted in a dendrogram, which we divided into 21 groups (Figure 5A–A' and Figure S4A–A'). The number of clusters was chosen by visual inspection so that most groups contained one or a few cell types based both directly on morphological stereotypy and on our reading of the previous literature (Otsuna and Ito, 2006). One of the main differences between the groups was the central brain neuropils they innervated. For example, groups I to III were the only ones that innervated the anterior optic tubercle (AOTU) (Figure 5A', B). We further investigated these groups to determine if we could match them to previously reported VPN classes (Otsuna and Ito, 2006) and to determine if central brain innervation was a major differentiating characteristic between classes.

**Lobula-, AOTU- and PVLP-innervating uVPns** We took neuron skeletons from groups I–III and isolated just the uVPN axon arbors that co-localized with central brain neuropils: the AOTU and posterior ventrolateral protocerebrum (PVLP) (Figure 5B). We then carried out a new clustering based on all-by-all NBLAST scores of these partial skeletons, cutting the dendrogram at a level at which we could obtain isolated types (groups 1–7). A clear distinction between neurons that innervated the PVLP (groups 1, 2) and those that extensively innervated the AOTU (groups 3–6) was evident. Table 1 lists the full correspondences between cluster groups 1–7 and the types previously described by the comprehensive analysis of Otsuna and Ito (2006). Our analysis divides the LC10 uVPN class into 5

Dendrogram group	Neuron type	Comments
1	LC6	Lobula innervation, dorsal cell bodies, axons follow the anterior optic tract to the AOTU, turning ventrally midway to innervate the lateral PVLP.
2	LC9	As group 1, but terminating in the medial PVLP rather than extending laterally.
3	LC10B (possibly different subtype to 4)	Dorsal AOTU innervation, ventral to group 4.
4	LC10B (possibly different subtype to 3)	Dorsal AOTU innervation, dorsal to group 3.
5	New LC10 subclass	The most ventral AOTU innervation. Similar to group 7, but ventral to it in the AOTU.
6	LC10A	Axons project through ventral AOTU, turn sharply dorsally to terminate in the dorsal AOTU.
7	New LC10 subclass	Ventral AOTU innervation, dorsal to group 5.

Table 1: Correspondences between hierarchical clustering groups of AOTU- and PVLP-innervating uVPNs via NBLAST scores and previously determined neuron types

subgroups: groups 5 and 7 define new LC10 subclasses that have not previously been reported, while groups 3 and 4 split the previously reported LC10B subclass into two distinct types (see also Figure 5C). Thus our computational analysis is able to identify new cell types even for intensively studied neuronal classes.

**Lobula-, PVLP- and PLP-innervating uVPNs** We performed a similar analysis with uVPN groups that had dendritic innervation restricted to the lobula and axons projecting to the PVLP and posterior lateral protocerebrum (PLP) (groups IV, VI, VII and XI) (Figure S4A'). Following the same strategy, we re-clustered neurons based on NBLAST scores calculated only for the axon arbors that overlapped with the PVLP or PLP (Figure S4B). A dendrogram divided into seven groups (A–G) resulted in clearly distinct neuron morphologies, with neurons innervating different glomeruli in the PVLP (Figure S4B–C). Table 2 details the full correspondences between cluster groups A–F and previously described neuron types (Otsuna and Ito, 2006). Neurons in group F define a new type that resembles but is distinct from LT12 neurons: they project along the PVLP posteriorly to LT12 and terminate in a lateral region of the superior posterior slope (SPS).

**Bilateral VPns** In addition to the analysis of the uVPNs, we also performed a hierarchical clustering of the bilateral VPns (Figure 5A''). We divided the resulting dendrogram in 8 groups (A–H), which allowed us to obtain unique neuron types, based on the cell morphology. We were able to match group

Dendrogram group	Neuron type	Comments
A	LC12 (possibly different subtype to B)	Innervation in most lateral and anterior PVLP glomerulus.
B	LC12 (possibly different subtype to A)	Innervation in more posterior and medial PVLP glomerulus than group A.
C	LC4	Innervates a more medial PVLP glomerulus than LC12.
D	LT12	Tentative match. Class was identified in <a href="#">Otsuna and Ito (2006)</a> based on a single neuron.
E	LC11	Innervates more dorsal PVLP glomerulus than LC12. Extends along the posterior PVLP, with a sharp anterior turn. Terminates with a blunt-stick like ending in the lateral PVLP.
F	New LT subclass	Similar to LT12, but with projections posterior to it, terminating in the lateral region of the superior posterior slope.
G	Unmatched	Do not correspond to a single type.

Table 2: Correspondences between hierarchical clustering groups of PVLP- and PLP-innervating uVPNs via NBLAST scores and previously determined neuron types

B to LC14 neurons ([Otsuna and Ito, 2006](#)). These neurons connect the two lobulas, have dorsal cell bodies and arborize on the surface of the ipsilateral lobula, projecting to the contralateral one via the great commissure.

**VPN summary** Our analysis of VPN neurons has demonstrated that the similarity algorithm we have developed has multiple applications. In addition to searches against a whole neuron, similarity searches performed with only part of a neuron are useful to highlight morphological features that might be most important for defining neuron classes. For VPNs, we used neuropil innervation to define groups of similar neurons. We were able to match 11 of these groups to known VPN types, and furthermore described two new subclasses and four subtypes of uVPNs, showing that this type of analysis can also reveal unknown neuron classes.

### Auditory neurons

Auditory projection neurons (PNs) are characterized by their innervation of the primary or secondary auditory neuropils, the antennal mechanosensory and motor center (AMMC) and the inferior ventrolateral protocerebrum (IVLP or wedge) and several distinct types have been described based on anatomical and physiological features ([Yorozu et al., 2009](#); [Lai et al., 2012](#); [Kamikouchi et al., 2006, 2009](#)). Similarly to the VPNs, we wanted to explore if we would be able to identify known and new auditory PNs using NBLAST searches. In this case, we employed a two search strategy. First, for 5 of the auditory types, we used the FlyCircuit neuron named by [Lai et al. \(2012\)](#) as the seed neuron for the first search. Candidate neurons were selected using strict anatomical criteria. A second search was then done using these candidates as query neurons and collecting all high scorers (score over 0.5). For each of the 5 auditory types, hierarchical clustering of the hits revealed new subtypes of known auditory PN types that differed mainly in the regions innervated in either the AMMC or IVLP (Figure S5). For the AMMC-IVLP PN2 type, we identified 6 possible distinct subtypes, with only one of these matching the seed neuron (Figure S5E). Table S1 details the full correspondences between cluster groups and previously described auditory PN types ([Lai et al., 2012](#)).

### mAL neurons

The *fruitless*-expressing mAL neurons are sexually dimorphic interneurons that are known to regulate wing extension by males during courtship song ([Koganezawa et al., 2010](#); [Kimura et al., 2005](#)). Males have about 30 neurons, but there are only 5 in females. Although the gross neuronal morphology is similar in both sexes, both axonal and dendritic arborisations are located in distinct regions, likely altering input and output connectivity. For example it appears that only male neurons receive gustatory input via dendritic arbors in the subesophageal ganglion ([Koganezawa et al., 2010](#)). We investigated whether clustering could distinguish male and female neurons and identify male subtypes. From a search with a seed mAL neuron, *fru-M-500159*, we collected 41 hits with a mean NBLAST score greater than 0.2. Hierarchical clustering divided these into 2 groups, cleanly separating male from female neurons (Figure 6A–B). In order to explore whether there were distinct types of male neuron, we carried out clustering using partial skeletons containing only the axonal and dendritic arbors, omitting the primary neurites and axon that are common to all neurons (Figure 6C). We identified 3 main types (groups I–III) and 2 subtypes of male mAL neurons. The 3 types differed in the length of the ipsilateral



ventral projection; this has feature previously been proposed as the basis of a qualitative classification of mAL neurons (Kimura et al., 2005). However all types also showed reproducible differences in the exact location of their axon terminal arbors, a feature that also distinguished subtypes IA and IB. Our analysis therefore suggests that the population of male neurons includes types with correlated differences in input-output connectivity.

## P1 neurons

P1 neurons are the most significantly dimorphic *fruitless*-expressing neurons. Male P1 neurons are involved in the initiation of male courtship behavior while female P1 neurons degenerate during development due to the action of *doublesex* (Kimura et al., 2008). There are around 20 P1 neurons that develop from the pMP-e *fruitless* neuroblast clone, exhibiting its distinctive primary neurite, that runs dorsally and then bends anteroventrally. They have extensive bilateral arborizations in the PVLP and ring neuropil, partially overlapping with the male-specific enlarged brain regions (MER) (Kimura et al., 2008; Cachero et al., 2010) (Figure 6D). Direct activation of P1 neurons shows that they play a critical role in initiating male courtship (Kohatsu et al., 2011; von Philipsborn et al., 2011; Inagaki et al., 2014; Bath et al., 2014), but until now they have been treated as homogeneous neuronal class. We therefore investigated whether clustering could identify anatomical subtypes. We first identified P1 neurons by searching the FlyCircuit dataset with a tracing of the distinctive primary neurite of a pMP-e clone (Cachero et al., 2010). Hierarchical clustering of the top hits identified a subset consisting solely of P1 neurons. Hierarchical clustering of these P1 neurons, cutting the dendrogram at height 1 produced 10 groups (Figure 6D'). Nine of these (1–9) contain only male *fru*-GAL4 neurons as expected. The 9 male groups have the same distinctive primary neurite and send contralateral axonal projections through the arch (Yu et al., 2010b) with extensive arborisations in the strongly male enlarged regions of the brain (Cachero et al., 2010). However each group shows a highly distinctive pattern of dendritic and axonal arborisations suggesting that they are likely to integrate distinct sensory inputs and to connect with distinct downstream targets. It will therefore be very interesting to test whether these anatomical classes have distinct functional roles in male sexual behavior.

Intriguingly group 10, consists only of female neurons, including two female *fru*-GAL4 and two other drivers. Their morphology is clearly similar to but distinct from group 9 neurons, suggesting that a small population of neurons that share anatomical (and likely developmental) features with male P1 neurons is also present in females.

## Matching GAL4 traces to single neurons

We have demonstrated NBLAST queries using a single neuron, group of neurons or a traced neurite. However, another important use case is to identify neuronal types labelled by genetic driver lines that may contain multiple neuronal classes and where the detailed morphology of individual neurons cannot be ascertained. We have developed a simple and rapid approach, which we exemplify with the GAL4 line R18C12 (Jenett et al., 2012). The expression pattern includes an obvious bilateral dorsal tract associated with a specific cluster of neurons (Figure 7A). We carried used Vaa3D (Peng et al., 2014) to trace this tract from a registered confocal stack downloaded from [virtualflybrain.org](http://virtualflybrain.org) and transformed it into the FCWB space. An NBLAST search of this tracing identified three very similar FlyCircuit neurons, which completely overlapped with the expression of R18C12. These three neurons appear to

correspond to different subtypes that vary in their terminal arborizations (Figure 7A–B). Peng et al. (2014) have recently published a projectome of 9198 such tracings from 1107 Gal4 lines, each of which can be used to retrieve (in seconds) matching neurons from the FlyCircuit dataset (data not shown).

## Superclusters and Exemplars to Organize Huge Data

In the previous examples we have shown that NBLAST is a powerful tool to identify known and uncover new neuron types when analyzing specific neuron superclasses within large datasets. However, subsetting the dataset in order to isolate the chosen neurons required directed searches with a seed neuron or queries for neuropil overlay, requiring considerable time. We wished to establish a method that would allow us to organize large datasets, extracting the main types automatically, and retain information on the similarity between types and subtypes, while providing a quicker way to analyze large number of images. We used the affinity propagation method of clustering (Frey and Dueck, 2007), combined with hierarchical clustering to achieve this. Affinity propagation is an unsupervised clustering method which, given a similarity measure, collapses groups of many alike objects and, for each group, finds a single exemplar representative of all the other members (Figure 8A). Applying this method to the 16,129 neurons in the FlyCircuit dataset resulted in 1,052 clusters. These differed in size and similarity, but each contained on average 10 neurons and a similarity of 0.56 (Figure 8B). We then performed a hierarchical clustering of the 1,052 exemplars, and divided the dendrogram into three groups (A–C), in order to obtain large subdivisions which we could analyze in more detail (Figure 8C). Apart from a few stray neurons, group A contained all the optic lobe and visual projection neurons, group C the neurons in the protocerebrum and group B the remaining neurons in the central brain. Hierarchical clustering of central brain exemplars in groups B and C revealed large superclasses of neuron types when we divided the dendrogram in 14 groups (I–XIV) (Figure 8D–D'). These superclasses included, for example, central complex neurons (I), P1 neurons (II), 2 groups of KCs ( $\gamma$  and  $\alpha'/\beta'$  and  $\alpha/\beta$ ) (IV–V), auditory neurons (VIII), and uPNs (XI).

The affinity propagation clusters are useful for revealing differences in morphology and possible subtypes between neurons in the same class (Figure 8E). Three clusters corresponding to 82 auditory neurons of the type antennal mechanosensory and motor center AMMC-IVLP projection neuron 1 (AMMC-IVLP PN1) (Lai et al., 2012) innervate different regions of the AMMC, and one of the clusters (red) does not extend as far laterally in the inferior ventrolateral protocerebrum or wedge (IVLP) as the other two. These disparities in morphology in the AMMC-IVLP PN1 have not been described previously. In the case of uVPNs, the class LC10B that innervates the dorsal AOTU corresponds to 11 clusters and 121 neurons, and there are clear differences in the regions innervated in the lobula, from dorsolateral to anteriomedial (Figure 5B, 8E). For the class LC4 that mostly innervates the PVLP, we found 11 exemplars, corresponding to 98 neurons (Figure S4B, 8E). Although a few members of these clusters are not LC4, the lamination in the lobula is evident, from dorsal to ventral.

We have shown that combining affinity propagation with hierarchical clustering is an effective way to organize and explore large datasets, by condensing information into a single exemplar and by retaining the ability to move up or down in the hierarchical tree, allowing the analysis of superclasses or more detailed subtypes.

## Discussion

The challenge of mapping and cataloguing the full spectrum of neuronal types in the brain depends not only on the ability to recognize similar neurons, by shape and position, but also on establishing methods that facilitate unbiased identification of neuron types from pools of thousands or millions of individual neurons. Comparison of 3D structures like neurons is a much harder problem than alignments of linear protein or DNA. The spatial position of a neuron within the brain is an essential determinant of its function and synaptic partners. For this reason, comparison of neurons must include not only morphology, but also position within the brain. Our approach included a pre-processing step that transformed all neurons into a common 3D reference space, using non-rigid image registration. Although this can be time and resource intensive, our results demonstrate that this is an effective way to achieve accurate comparisons. A neuron search algorithm should therefore be: (1) accurate, with hits being biologically meaningful; (2) fast and computationally inexpensive to allow multiple searches; (3) able to provide the user with an interactive search method and (4) generally applicable. Here, we have described NBLAST, a neuron search algorithm that satisfies all these criteria.

First, the algorithm correctly distinguishes closely related subtypes across a range of major neuron groups, with an accuracy of 97.6% in the case of olfactory projection neurons. Unsupervised clustering of these neurons, based on NBLAST scores, correctly organized neurons into described types. We did find, however, that the size of a neuron influences the accuracy of algorithm, especially for smaller neurons, even when using the normalized score. One future research area will be to convert the raw scores that we have used into an expectation (E) value (cf. BLAST), that would directly account for the size of a neuron.

Second, NBLAST searches are very fast, with pairwise comparisons taking about 2ms on a laptop computer, with queries against the whole 16,000 neuron dataset taking about 30s. Furthermore for defined datasets all-by-all scores can be pre-computed allowing immediate retrieval of NBLAST scores for highly interactive analysis. Large scale efforts to map the fly brain, combining sparse genetic labeling with light microscopy are gathering pace with over 20 Tb of data publicly available. This could be similar to the rapid growth experience in the early days of protein databases such as UniProtKB (Wu et al., 2006; Consortium, 2014b,a). As neuron databases get bigger, so will the time required to query a single neuron against the full set. While this increase will be linear, the storage space required to save pre-computed scores will increase quadratically. Currently, our all-by-all score matrix for 16,000 neurons requires ~2 GB of disk space. One effective approach to handle this will be to compute sparse similarity matrices for example storing on only the top  $n$  hits for a given neuron, an approach often taken for genome-wide precomputed BLAST scores. Alternatively, queries could be computed only against the non-redundant set of neurons that collectively embody the structure of the brain, similarly to the strategy employed by UniProt (Suzek et al., 2007). At most, this set could not exceed 50,000 neurons (due to the strong bilateral symmetry in the fly brain) and we expect that it would in practice not need to exceed 5,000. Our clustering of all ~16,000 neurons of the FlyCircuit dataset identified ~1,000 exemplars providing a non-redundant data set that could be used for rapid searches.

Third, our method permits a variety of different search strategies accuracy. We have demonstrated that NBLAST searches can use a whole neuron or only part of a neuron as a query. The latter approach can use either highly conserved features such as axon tracts (e.g. for identifying candidate neurons from a Gal4 line Figure 7) or to distinguish closely related neuronal types by their terminal arbors (Figure 6).

Finally, one important question is obviously the extent to which our approach can be generalized. This issue largely reduces to the relationship between the length scales of neurons being examined and their absolute spatial stereotypy. Our method implicitly assumes spatial co-localization of related neurons; this is enforced in our input data by the use of image registration. Our search strategy should therefore be appropriate for any situation in which neuronal organization is highly stereotyped at the length scale of the neuron under consideration. There is already strong evidence that this is true across large parts of the brain for simple vertebrate models like the larval zebrafish: indeed [Portugues et al. \(2014\)](#) used exactly the same registration software that we have used in flies to demonstrate highly spatially stereotyped visuomotor activity patterns. Preliminary analysis (GSXEJ and JDM, <https://github.com/jefferis/nat.examples/tree/master/05-miyasaka2014>) suggests that our method can be applied directly to olfactory projectome data ([Miyasaka et al., 2014](#)) from larval zebrafish. Mouse gene expression ([Lein et al., 2007](#)) and long range connectivity also show global spatial stereotypy as evidenced by recent atlas studies combining sparse labeling and image registration ([Zingg et al., 2014](#); [Susaki et al., 2014](#); [Oh et al., 2014](#)). Our method should allow simple querying and hierarchical organization of these datasets with relatively little modification beyond calculating an appropriate scoring matrix.

However there are clearly situations in which global brain registration is not an appropriate starting point. For example the vertebrate retina has both a laminar and a tangential organization. Recently [Sümbül et al. \(2014\)](#) have introduced a registration strategy that demonstrates that the lamination of retinal ganglion cells in mouse retina is spatially stereotyped to the nearest micron – an even higher degree of spatial stereotypy than has been reported in flies after global brain registration ([Jefferis et al., 2007](#)). However retinal interneurons and ganglion cells are organized in mosaics across the retinal surface (typically referred to as the XY plane). Therefore global registration is not appropriate in this axis, rather it is necessary to align neurons into a virtual column – for example by the simple expedient of centering neurons on the XY position of their cell body or the centroid of all their segments.

The situation is similar for *Drosophila* columnar neurons in the outer neuropils of optic lobe, for which our method will need modification for optimal results. Inputs and some aspects of downstream processing are organized into about 800 parallel columns (reviewed in [Paulk et al., 2013](#)). Two neurons of exactly the same type may not be co-localized if they are labelled in different columns. There are two ways that we envisage this situation can be handled. The first would be to carry out a local re-registration, that maps each column onto a single canonical column. The second would be to amass sufficient data that neurons from neighboring columns would tile the brain, enabling their identification as a related group by standard clustering or graph theoretic approaches.

The aim of cataloguing all neuron types in the brain relies not only on an accurate algorithm to find similar neurons, but also on having an easy and unbiased method to distinguish neuron types and/or subtypes. This is a challenging problem, but morphological approaches may eventually provide unambiguous automated classification. [Sümbül et al. \(2014\)](#) recently explored the issue of defining the optimal cut height for morphological clustering of mouse retinal ganglion cells establishing a reliable approach for these specific neurons. We have demonstrated the applicability of NBLAST across a very wide range of neuronal classes in identifying known and novel neuronal types using hierarchical clustering. We could identify candidate types by cutting the dendrogram at a specified height. This process of identification relied on extensive data exploration and iteration. We would not expect to find a unique value for the dendrogram height at which neuron types would always be identified due to the differences in the extent of morphological variability within a type. Even so, for the examples

we present, except for Kenyon cells, we found a range of heights between 0.7 and 2. For olfactory projection neurons, where types can be unambiguously identified by their uniglomerular innervation, the dendrogram cut height was 0.75. The range of values we found will guide future exploration for other neuron types and datasets, although this process will still require iterative analysis and manual verification.

## Experimental Procedures

### Image Preprocessing

The flycircuit.tw team supplied 16,226 raw confocal stacks in the Zeiss LSM format on a single 2 TB hard drive in April 2011. Each LSM stack was first uncompressed, then read into Fiji/ImageJ (<http://fiji.sc/Fiji>) where the channels were split and resaved as individual gzip compressed NRRD (<http://teem.sourceforge.net/nrrd>) files. Where calibration information was missing from the LSM file metadata, we used a voxel size of (0.318427, 0.318427, 1.00935) microns as recommended by the FlyCircuit team. There were two important issues to solve before images could be used for registration: 1) identifying which image channel contained the anti-Dlg (discs large 1) counterstaining revealing overall brain structure and 2) determining whether the brains had been imaged from anterior to posterior, or posterior to anterior. The first issue could be solved by exporting the metadata associated with each LSM file using the LOCI bioformats (<http://loci.wisc.edu/software/bio-formats>) plugin for Fiji and developing some heuristics to automate the identification of the channel sequence; for a minority of images this metadata was missing and the channel order was determined manually. The second issue, slice order, could not be determined automatically from the image metadata. We therefore made maximum intensity projections (using the unu tool, <http://teem.sourceforge.net/unrrdu>) along the Z axis of the channel corresponding to the labelled neuron for each stack. Each projection was then compared with the matching thumbnail available from the flycircuit.tw website. The correlation score between the projection and thumbnail images was calculated both with and without a mirror flip across the YZ plane; a large correlation score for only one orientation was used as evidence for a given slice ordering. A small number of ambiguous results were verified manually. We successfully preprocessed 16,204/16,226 total images i.e. a 0.14 % failure rate. 12 failures were due to mismatches that could not be resolved between the segmented neuron present in the LSM file and the thumbnail image for the neuron identifier on the flycircuit.tw website; the remaining 10 failures were due to physical offsets between the brain and GFP channels or corrupt image data.

### Template Brain

The template brain (FCWB) was constructed by screening for whole brains within the FlyCircuit dataset, and manually selecting a pool of brains that appeared of good quality when the stacks were inspected. Separate average female and average male template brains were constructed from 17 and 9 brains, respectively using the CMTK (<http://www.nitrc.org/projects/cmtk>) `avg_adm` function which takes a single brain as a seed. After five iterations the resultant average male and average female brains were placed in an affine symmetric position within their image stacks so that a simple horizontal ( $x$ -axis) flip of either template brain resulted in an almost perfect overlap of



left and right hemispheres. Finally the two sex-specific template brains were then averaged (with equal weight) to make an intersex template brain using the same procedure. Since the purpose of this template was to provide an optimal registration target for the [flycircuit.tw](http://flycircuit.tw) dataset, no attempt was made to correct for the obvious disparity between the XY and Z voxel dimensions common to all the images in the dataset. The scripts used for the construction of the template are available at <https://github.com/jefferislab/MakeAverageBrain>.

## Image Registration

Image registration of the Dlg neuropil staining used a fully automatic intensity-based (landmark free) 3D image registration implemented in the CMTK toolkit available at <http://www.nitrc.org/projects/cmtk> (Rohlfing and Maurer, 2003; Jefferis et al., 2007). An initial linear registration with 9 degrees of freedom (translation, rotation and scaling of each axis) was followed by a non-rigid registration that allows different brain regions to move somewhat independently, subject to a smoothness penalty (Rueckert et al., 1999). It is our experience that obtaining a satisfactory initial linear registration is crucial. All registrations were therefore manually checked by comparing the reformatted brain with the template in Amira (academic version, Zuse Institute, Berlin), using ResultViewer <https://bitbucket.org/jefferis/resultviewer>. This identified about 10 % of brains which did not register satisfactorily. For these images a new affine registration was calculated using a Global Hough Transform function available in Amira; the result of this affine transform was again manually inspected. In the minority of cases where this approach failed, a surface based alignment was calculated in Amira after manually aligning the two brains. Once a satisfactory initial affine registration was obtained, a non-rigid registration was calculated for all brains. Finally each registration was checked manually in Amira against the template brain. The result of this sequential procedure was that we successfully registered 16,129/16,204 preprocessed images, giving a registration failure rate of 0.46 %.

## Image Postprocessing

The confocal stack for each neuron available at <http://flycircuit.tw> includes an 8 bit image containing a single (semiautomatically) segmented neuron prepared by Chiang et al. (2011). This image was downsampled by a factor of 2 in  $x$  and  $y$ , binarized with a threshold of 1 and then skeletonized using the Fiji plugin ‘Skeletonize (2D/3D)’. Dot properties for each neuron skeleton were extracted following the method in Masse et al. (2012), using the dotprops function of our new [nat](#) package for R. This converted each skeleton into segments, described by its location and tangent vector. Neurons on the right side of the brain were flipped to left by applying a mirroring and a flipping registration as described in Manton et al. (2014). The decision of whether to flip a neuron depended on earlier assignment of each neuron to a brain hemisphere using a combination of automated and manual approaches. Neurons whose cell bodies were more than 20  $\mu\text{m}$  away from the mid-sagittal YZ plane were automatically defined as belonging to the left or right hemisphere. Neurons whose cell bodies were inside this 40  $\mu\text{m}$  central corridor were manually assigned to the left or right sides, based on the position of the cell body (right or left side), path taken by the primary neurite, location and length of first branching neurite. For example, neurons that had a cell body on the midline with significant innervation from the first branching neurite near the cell body on the left hemisphere, with the rest of the arborisation on the right, were assigned to the left side and not flipped. On the other hand, neurons



with similar morphology to these but in which the first branching neurite is small, compared to the total innervation, was assigned to the right and flipped. The cell body positions used were based on those published on the <http://flycircuit.tw> website for each neuron; these positions are in the space of the FlyCircuit female and male template brains (typical\_brain\_female and typical\_brain\_male). In order to transform them into the FCWB template that we constructed, affine bridging registrations were constructed from the typical\_brain\_female and typical\_brain\_male brains to FCWB and the cell body positions were then transformed to this new space. Since these cell body positions depend on two affine registrations (one conducted by Chiang et al. (2011) to register each sample brain onto either their typical\_brain\_female or typical\_brain\_male templates and a second carried out by us to map those template brains onto our FCWB template) these positions are likely accurate only to  $\pm 5$  microns in each axis.

## Neuron Search

The neuron search algorithm is described in detail in Results and Figure 1. The reference implementation that we have written is the `nblast` function in the R package `nat.nblast`, which depends on our `nat` package (Jefferis and Manton, 2014). Fast nearest neighbor search, an essential primitive for the algorithm uses the RANN package (Jefferis, 2014), a wrapper for the Approximate Nearest Neighbor (ANN) C++ library (Mount and Arya, 2006). The scoring matrix that we used for FlyCircuit neurons was constructed by taking 150 DL2 projection neurons, which define a neuron type at the finest level, and calculating the joint histogram of distance and absolute dot product for the  $150 \times 149$  combinations of neurons, resulting in  $1.4 \times 10^7$  measurement pairs; the number of counts in the histogram was then normalized (i.e. dividing by  $1.4 \times 10^7$ ) to give a probability density,  $p_{\text{match}}$ . We then carried out a similar procedure for 5,000 random pairs of neurons sampled from the FlyCircuit dataset to give  $p_{\text{rand}}$ . Finally the scoring matrix was calculated as  $\log_2 \frac{p_{\text{match}} + \epsilon}{p_{\text{rand}} + \epsilon}$  where  $\epsilon$  (a pseudocount to avoid infinite values) was set to  $1 \times 10^{-6}$ .

## Clustering

We employed two different methods for clustering based on normalized NBLAST scores. We used Ward's method for hierarchical clustering, using the default implementation in the R function `hclust`. This method minimizes the total within-cluster variance, and at each step the pair of entities or clusters with the minimum distance between clusters are merged (Ward Jr, 1963). The resulting dendrograms were cut at a single selected height chosen for each case to separate neuron types or subtypes. This value is shown as a dashed line in all dendrograms. By default, R plots the square of the Euclidean distance as the  $y$  axis, but in the plots shown, the height of the dendrogram corresponds to the unsquared distance.

For the analysis of the whole dataset, we used the affinity propagation method. This is an iterative method which finds exemplars which are representative members of each cluster and does not require any *a priori* input on the final number of clusters (Frey and Dueck, 2007) as implemented in the R package `apcluster` (Bodenhofer et al., 2011). The input preference parameter ( $p$ ) can be set before running the clustering. This parameter reflects the tendency of data samples to become an exemplar, and affects the final number of clusters. In our analysis, we used  $p = 0$ , since this is the value where on average matched segments are equally likely to have come from matching and non-matching neurons. Empirically this parameter produced clusters that, for the most part, grouped neurons of the same type

according to biological expert opinion.

## Neuron Tracing

Neuron tracing was carried out in Amira (commercial version, FEI Visualization Sciences Group, Merignac, France) using the hxskeletonize plugin (Evers et al., 2005) or in Vaa3D (Peng et al., 2014) on previously registered image data. Traces were then loaded into R using the nat package. When necessary, they were transformed into the space of the FCWB template brain using the approach of Manton et al. (2014).

## Computer Code and Data

The image processing pipeline and analysis code used two custom packages for the R statistical environment (<http://www.r-project.org>) <https://github.com/jefferis/nat> and <https://github.com/jefferis/nat.as> that coordinated processing by the low level registration (CMTK) and image processing (Fiji, unu) software mentioned above. NBLAST neuron search is implemented in a third R package available at <https://github.com/jefferislab/nat.nblast>. Analysis code specific to the flycircuit dataset is available in a dedicated R package <https://github.com/jefferis/flycircuit>, with a package vignette showcasing the main tools that we have developed. Further details of these supplemental software and the associated data are presented at <http://jefferislab.org/si/nblast>. The registered image dataset can be viewed in the stack viewer of the <http://virtualflybrain.org> website and all 16,129 registered single neuron images will be available at <https://jefferislab.org/si/nblast> or on request to GSXEJ on a hard drive; the unregistered data remain available at <http://flycircuit.tw>.

## Acknowledgments

We first of all acknowledge the flycircuit.tw team for generously providing the raw image data associated with Chiang et al. (2011). Images from FlyCircuit were obtained from the NCHC (National Center for High-performance Computing) and NTHU (National Tsing Hua University), Hsinchu, Taiwan. We thank %XYZ and members of the Jefferis lab for comments on the manuscript, Jake Grimmett and Toby Darling for assistance with the LMB's compute cluster and Torsten Rohlfing for discussions about image analysis and registration. We thank the Virtual Fly Brain project for their help in linking and incorporating some of the results of this study in the <http://virtualflybrain.org> website.

This study made use of the Computational Morphometry Toolkit, supported by the National Institute of Biomedical Imaging and Bioengineering. This work was supported by the Medical Research Council [MRC file reference U105188491] and a European Research Council Starting Investigator Grant to GSXEJ, who is an EMBO Young Investigator.

## References

- Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local alignment search tool. *Journal of Molecular Biology* 215, 403–410.
- Ascoli, G.A., Donohue, D.E., and Halavi, M. (2007). NeuroMorpho.Org: a central resource for neuronal morphologies. *J Neurosci* 27, 9247–51.
- Aso, Y., Grübel, K., Busch, S., Friedrich, A.B., Siwanowicz, I., and Tanimoto, H. (2009). The mushroom body of adult *Drosophila* characterized by GAL4 drivers. *Journal of neurogenetics* 23, 156–172.
- Badea, T.C., Wang, Y., and Nathans, J. (2003). A noninvasive genetic/pharmacologic strategy for visualizing cell morphology and clonal relationships in the mouse. *J Neurosci* 23, 2314–22.
- Badea, T.C., and Nathans, J. (2004). Quantitative analysis of neuronal morphologies in the mouse retina visualized by using a genetically directed reporter. *Journal of Comparative Neurology* 480, 331–351.
- Basu, S., Condrón, B., and Acton, S.T. (2011). Path2Path: hierarchical Path-Based analysis for neuron matching. In *Biomedical Imaging: From Nano to Macro, 2011 IEEE International Symposium on (IEEE)*, pp. 996–999.
- Bath, D.E., Stowers, J.R., Hörmann, D., Poehlmann, A., Dickson, B.J., and Straw, A.D. (2014). Fly-MAD: rapid thermogenetic control of neuronal activity in freely walking *Drosophila*. *Nat Methods* 11, 756–62.
- Bodenhofer, U., Kothmeier, A., and Hochreiter, S. (2011). APCluster: an R package for affinity propagation clustering. *Bioinformatics* 27, 2463–2464.
- Bota, M., and Swanson, L.W. (2007). The neuron classification problem. *Brain research reviews* 56, 79–88.
- Brown, K.M., Barrionuevo, G., Canty, A.J., De Paola, V., Hirsch, J.A., Jefferis, G.S.X.E., Lu, J., Snippe, M., Sugihara, I., and Ascoli, G.A. (2011). The DIADEM Data Sets: Representative Light Microscopy Images of Neuronal Morphology to Advance Automation of Digital Reconstructions. *Neuroinformatics* .
- Cachero, S., Ostrovsky, A.D., Yu, J.Y., Dickson, B.J., and Jefferis, G.S.X.E. (2010). Sexual dimorphism in the fly brain. *Curr Biol* 20, 1589–601.
- Cajal, S.R., and Azoulay y, L. (1911). *Histologie du système nerveux de l’homme et des vertébrés* (A. Maloine).
- Cardona, A., Saalfeld, S., Arganda, I., Pereanu, W., Schindelin, J., and Hartenstein, V. (2010). Identifying neuronal lineages of *Drosophila* by sequence analysis of axon tracts. *J Neurosci* 30, 7538–53.
- Caron, S.J.C., Ruta, V., Abbott, L.F., and Axel, R. (2013). Random convergence of olfactory inputs in the *Drosophila* mushroom body. *Nature* 497, 113–7.

- Chiang, A.S., Lin, C.Y., Chuang, C.C., Chang, H.M., Hsieh, C.H., Yeh, C.W., Shih, C.T., Wu, J.J., Wang, G.T., Chen, Y.C., Wu, C.C., Chen, G.Y., Ching, Y.T., Lee, P.C., Lin, C.Y., Lin, H.H., Wu, C.C., Hsu, H.W., Huang, Y.A., Chen, J.Y., Chiang, H.J., Lu, C.F., Ni, R.F., Yeh, C.Y., and Hwang, J.K. (2011). Three-dimensional reconstruction of brain-wide wiring networks in *Drosophila* at single-cell resolution. *Curr Biol* 21, 1–11.
- Consortium, U. (2014a). UniProtKB/Swiss-Prot protein knowledgebase release 2014\_07 statistics.
- Consortium, U. (2014b). UniProtKB/TrEMBL PROTEIN DATABASE RELEASE 2014\_07 STATISTICS.
- Conte, D., Foggia, P., Sansone, C., and Vento, M. (2004). Thirty years of graph matching in pattern recognition. *International journal of pattern recognition and artificial intelligence* 18, 265–298.
- Coombs, J., Van Der List, D., Wang, G.Y., and Chalupa, L. (2006). Morphological properties of mouse retinal ganglion cells. *Neuroscience* 140, 123–136.
- El Jundi, B., Heinze, S., Lenschow, C., Kurylas, A., Rohlfing, T., and Homberg, U. (2009). The locust standard brain: a 3D standard of the central complex as a platform for neural network analysis. *Frontiers in systems neuroscience* 3.
- Evers, J.F., Schmitt, S., Sibila, M., and Duch, C. (2005). Progress in functional neuroanatomy: precise automatic geometric reconstruction of neuronal morphology from confocal image stacks. *J Neurophysiol* 93, 2331–42.
- Fischbach, K.F., and Dittrich, A. (1989). The optic lobe of *Drosophila melanogaster*. I. A Golgi analysis of wild-type structure. *Cell and Tissue Research* 258, 441–475.
- Frey, B.J., and Dueck, D. (2007). Clustering by passing messages between data points. *science* 315, 972–976.
- Ganglberger, F., Schulze, F., Tirian, L., Novikov, A., Dickson, B., Bühler, K., and Langs, G. (2014). Structure-Based Neuron Retrieval Across *Drosophila* Brains. *Neuroinformatics* .
- Inagaki, H.K., Jung, Y., Hoopfer, E.D., Wong, A.M., Mishra, N., Lin, J.Y., Tsien, R.Y., and Anderson, D.J. (2014). Optogenetic control of *Drosophila* using a red-shifted channelrhodopsin reveals experience-dependent influences on courtship. *Nat Methods* 11, 325–32.
- Ito, K., Shinomiya, K., Ito, M., Armstrong, J.D., Boyan, G., Hartenstein, V., Harzsch, S., Heisenberg, M., Homberg, U., Jenett, A., et al. (2014). A systematic nomenclature for the insect brain. *Neuron* 81, 755–765.
- Jefferis, G.S.X.E., Vyas, R.M., Berdnik, D., Ramaekers, A., Stocker, R.F., Tanaka, N.K., Ito, K., and Luo, L. (2004). Developmental origin of wiring specificity in the olfactory system of *drosophila*. *Development* 131, 117–30.
- Jefferis, G. (2014). RANN k nearest neighbour search v2.3.0. Zenodo .
- Jefferis, G.S.X.E., and Manton, J.D. (2014). nat: NeuroAnatomy Toolbox R package. zenodo .

- Jefferis, G.S.X.E., Potter, C.J., Chan, A.M., Marin, E.C., Rohlfsing, T., Maurer, C.R.J., and Luo, L. (2007). Comprehensive maps of *Drosophila* higher olfactory centers: spatially segregated fruit and pheromone representation. *Cell* 128, 1187–1203.
- Jefferis, G.S.X.E., and Livet, J. (2012). Sparse and combinatorial neuron labelling. *Curr Opin Neurobiol* 22, 101–10.
- Jefferis, G.S., Marin, E.C., Stocker, R.F., and Luo, L. (2001). Target neuron prespecification in the olfactory map of *Drosophila*. *Nature* 414, 204–208.
- Jenett, A., Rubin, G.M., Ngo, T.T., Shepherd, D., Murphy, C., Dionne, H., Pfeiffer, B.D., Cavallaro, A., Hall, D., Jeter, J., et al. (2012). A GAL4-Driver Line Resource for *Drosophila* Neurobiology. *Cell reports* 2, 991–1001.
- Joseph, A.P., Shingate, P., Upadhyay, A.K., and Sowdhamini, R. (2014). 3PFDB+: improved search protocol and update for the identification of representatives of protein sequence domain families. *Database* 2014, bau026.
- Kahsai, L., and Zars, T. (2011). Learning and memory in *Drosophila*: behavior, genetics, and neural systems. *Int Rev Neurobiol* 99, 139–67.
- Kamikouchi, A., Inagaki, H.K., Effertz, T., Hendrich, O., Fiala, A., Göpfert, M.C., and Ito, K. (2009). The neural basis of *Drosophila* gravity-sensing and hearing. *Nature* 458, 165–171.
- Kamikouchi, A., Shimada, T., and Ito, K. (2006). Comprehensive classification of the auditory sensory projections in the brain of the fruit fly *Drosophila melanogaster*. *Journal of Comparative Neurology* 499, 317–356.
- Kepecs, A., and Fishell, G. (2014). Interneuron cell types are fit to function. *Nature* 505, 318–26.
- Kimura, K.I., Hachiya, T., Koganezawa, M., Tazawa, T., and Yamamoto, D. (2008). Fruitless and doublesex coordinate to generate male-specific neurons that can initiate courtship. *Neuron* 59, 759–769.
- Kimura, K.I., Ote, M., Tazawa, T., and Yamamoto, D. (2005). Fruitless specifies sexually dimorphic neural circuitry in the *Drosophila* brain. *Nature* 438, 229–233.
- Koganezawa, M., Haba, D., Matsuo, T., and Yamamoto, D. (2010). The Shaping of Male Courtship Posture by Lateralized Gustatory Inputs to Male-Specific Interneurons. *Current Biology* 20, 1–8.
- Kohatsu, S., Koganezawa, M., and Yamamoto, D. (2011). Female contact activates male-specific interneurons that trigger stereotypic courtship behavior in *Drosophila*. *Neuron* 69, 498–508.
- Kong, J.H., Fish, D.R., Rockhill, R.L., and Masland, R.H. (2005). Diversity of ganglion cells in the mouse retina: unsupervised morphological classification and its limits. *Journal of Comparative Neurology* 489, 293–310.
- Lai, J.S.Y., Lo, S.J., Dickson, B.J., and Chiang, A.S. (2012). Auditory circuit in the *Drosophila* brain. *Proc Natl Acad Sci U S A* 109, 2607–12.

- Le, Q., Ranzato, M., Monga, R., Devin, M., Chen, K., Corrado, G., Dean, J., and Ng, A. (2012). Building high-level features using large scale unsupervised learning. In International Conference in Machine Learning.
- Lee, T.C., Kashyap, R.L., and Chu, C.N. (1994). Building Skeleton Models via 3-D Medial Surface/Axis Thinning Algorithms. *CVGIP: Graph. Models Image Process.* 56, 462–478.
- Lee, T., Lee, A., and Luo, L. (1999). Development of the Drosophila mushroom bodies: sequential generation of three distinct types of neurons from a neuroblast. *Development* 126, 4065–4076.
- Lein, E.S., Hawrylycz, M.J., Ao, N., Ayres, M., Bensinger, A., Bernard, A., Boe, A.F., Boguski, M.S., Brockway, K.S., Byrnes, E.J., et al. (2007). Genome-wide atlas of gene expression in the adult mouse brain. *Nature* 445, 168–176.
- Lin, H.H., Lai, J.S.Y., Chin, A.L., Chen, Y.C., and Chiang, A.S. (2007). A map of olfactory representation in the Drosophila mushroom body. *Cell* 128, 1205–1217.
- Manning, G., Whyte, D.B., Martinez, R., Hunter, T., and Sudarsanam, S. (2002). The protein kinase complement of the human genome. *Science* 298, 1912–34.
- Manton, J.D., Ostrovsky, A.D., Goetz, L., Costa, M., Rohlfing, T., and Jefferis, G.S.X.E. (2014). Combining genome-scale Drosophila 3D neuroanatomical data by bridging template brains. *Bioarxiv preprint*.
- Marin, E.C., Jefferis, G.S., Komiyama, T., Zhu, H., and Luo, L. (2002). Representation of the Glomerular Olfactory Map in the Drosophila Brain. *Cell* 109, 243–255.
- Masse, N.Y., Turner, G.C., and Jefferis, G.S.X.E. (2009). Olfactory information processing in Drosophila. *Curr Biol* 19, R700–13.
- Masse, N.Y., Cachero, S., Ostrovsky, A., and Jefferis, G.S.X.E. (2012). A mutual information approach to automate identification of neuronal clusters in Drosophila brain images. *Frontiers in Neuroinformatics* 6.
- Mayerich, D., Bjornsson, C., Taylor, J., and Roysam, B. (2012). NetMets: software for quantifying and visualizing errors in biological network segmentation. *BMC Bioinformatics* 13 Suppl 8, S7.
- Migliore, M., and Shepherd, G.M. (2005). An integrated approach to classifying neuronal phenotypes. *Nature Reviews Neuroscience* 6, 810–818.
- Miyasaka, N., Arganda-Carreras, I., Wakisaka, N., Masuda, M., Sümbül, U., Seung, H.S., and Yoshihara, Y. (2014). Olfactory projectome in the zebrafish forebrain revealed by genetic single-neuron labelling. *Nat Commun* 5, 3639.
- Morante, J., and Desplan, C. (2008). The Color-Vision Circuit in the Medulla of Drosophila. *Current Biology* 18, 553–565.
- Mount, D.M., and Arya, S. (2006). ANN: A Library for Approximate Nearest Neighbor Searching. Version 1.1.1.

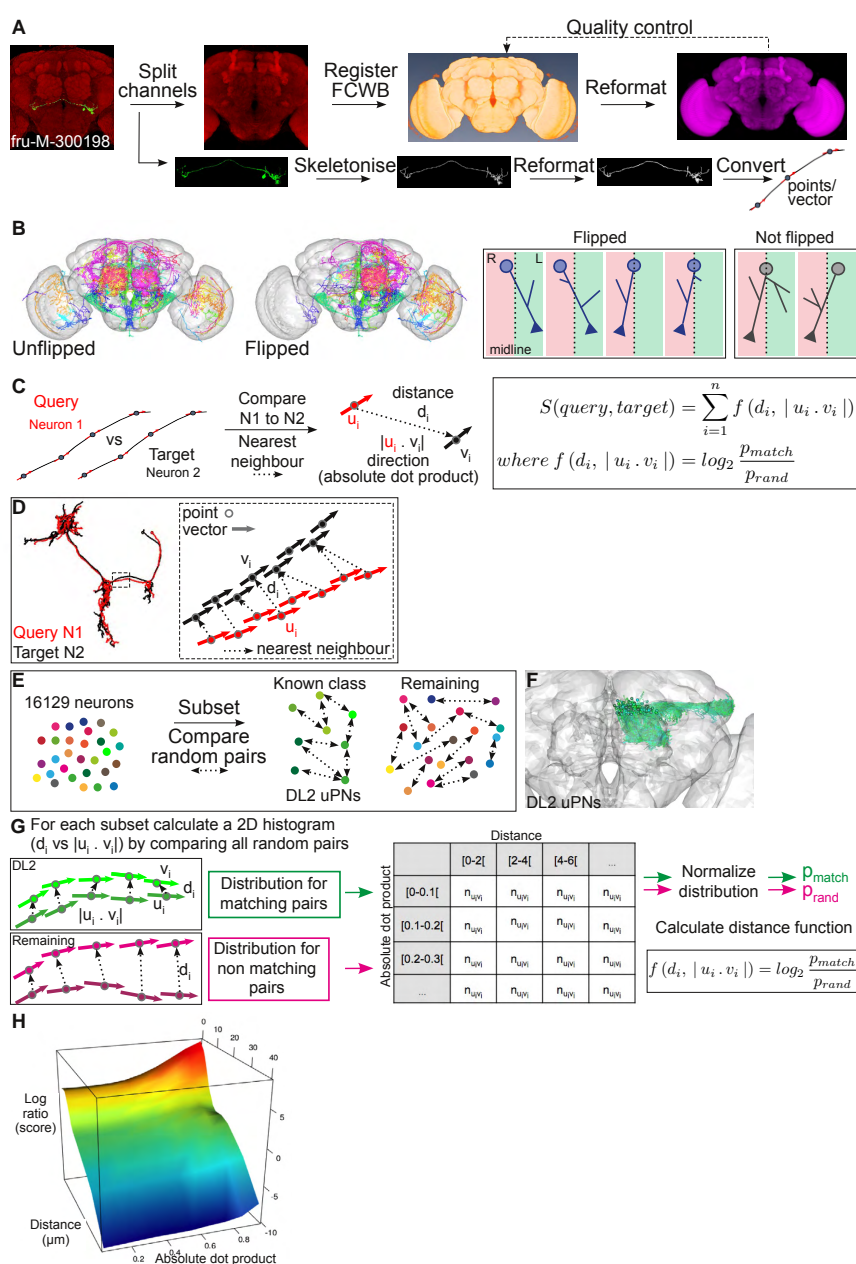


- Nelson, S.B., Sugino, K., and Hempel, C.M. (2006). The problem of neuronal cell types: a physiological genomics approach. *Trends Neurosci* 29, 339–45.
- Oh, S.W., Harris, J.A., Ng, L., Winslow, B., Cain, N., Mihalas, S., Wang, Q., Lau, C., Kuan, L., Henry, A.M., et al. (2014). A mesoscale connectome of the mouse brain. *Nature* 508, 207–214.
- Otsuna, H., and Ito, K. (2006). Systematic analysis of the visual projection neurons of *Drosophila melanogaster*. I. Lobula-specific pathways. *J Comp Neurol* 497, 928–958.
- Parekh, R., and Ascoli, G.A. (2013). Neuronal morphology goes digital: a research hub for cellular and system neuroscience. *Neuron* 77, 1017–38.
- Paulk, A., Millard, S.S., and van Swinderen, B. (2013). Vision in *Drosophila*: seeing the world through a model’s eyes. *Annual review of entomology* 58, 313–332.
- Pearson, W.R., and Lipman, D.J. (1988). Improved tools for biological sequence comparison. *Proc Natl Acad Sci U S A* 85, 2444–8.
- Peng, H., Tang, J., Xiao, H., Bria, A., Zhou, J., Butler, V., Zhou, Z., Gonzalez-Bellido, P.T., Oh, S.W., Chen, J., Mitra, A., Tsien, R.W., Zeng, H., Ascoli, G.A., Iannello, G., Hawrylycz, M., Myers, E., and Long, F. (2014). Virtual finger boosts three-dimensional imaging and microsurgery as well as terabyte volume image visualization and analysis. *Nat Commun* 5, 4342.
- Petilla Interneuron Nomenclature Group, Ascoli, G.A., Alonso-Nanclares, L., Anderson, S.A., Barionuevo, G., Benavides-Piccione, R., Burkhalter, A., Buzsáki, G., Cauli, B., Defelipe, J., Fairén, A., Feldmeyer, D., Fishell, G., Fregnac, Y., Freund, T.F., Gardner, D., Gardner, E.P., Goldberg, J.H., Helmstaedter, M., Hestrin, S., Karube, F., Kisvárdy, Z.F., Lambolez, B., Lewis, D.A., Marin, O., Markram, H., Muñoz, A., Packer, A., Petersen, C.C.H., Rockland, K.S., Rossier, J., Rudy, B., Somogyi, P., Staiger, J.F., Tamas, G., Thomson, A.M., Toledo-Rodriguez, M., Wang, Y., West, D.C., and Yuste, R. (2008). Petilla terminology: nomenclature of features of GABAergic interneurons of the cerebral cortex. *Nat Rev Neurosci* 9, 557–68.
- Portugues, R., Feierstein, C.E., Engert, F., and Orger, M.B. (2014). Whole-brain activity maps reveal stereotyped, distributed networks for visuomotor behavior. *Neuron* 81, 1328–1343.
- Rohlfing, T., and Maurer, C. R., J. (2003). Nonrigid image registration in shared-memory multiprocessor environments with application to brains, breasts, and bees. *IEEE Trans Inf Technol Biomed* 7, 16–25.
- Rowe, M., and Stone, J. (1976). Naming of neurones. Classification and naming of cat retinal ganglion cells. *Brain, behavior and evolution* 14, 185–216.
- Rueckert, D., Sonoda, L.I., Hayes, C., Hill, D.L., Leach, M.O., and Hawkes, D.J. (1999). Nonrigid registration using free-form deformations: application to breast MR images. *IEEE Trans Med Imaging* 18, 712–21.
- Rybak, J., Kuß, A., Lamecker, H., Zachow, S., Hege, H.C., Lienhard, M., Singer, J., Neubert, K., and Menzel, R. (2010). The digital bee brain: integrating and managing neurons in a common 3D reference system. *Frontiers in systems neuroscience* 4.

- Schindelin, J., Arganda-Carreras, I., Frise, E., Kaynig, V., Longair, M., Pietzsch, T., Preibisch, S., Rueden, C., Saalfeld, S., Schmid, B., Tinevez, J.Y., White, D.J., Hartenstein, V., Eliceiri, K., Tomancak, P., and Cardona, A. (2012). Fiji: an open-source platform for biological-image analysis. *Nat Methods* 9, 676–82.
- Sonnhammer, E.L., Eddy, S.R., Durbin, R., et al. (1997). Pfam: a comprehensive database of protein domain families based on seed alignments. *Proteins-Structure Function and Genetics* 28, 405–420.
- Sümbül, U., Song, S., McCulloch, K., Becker, M., Lin, B., Sanes, J.R., Masland, R.H., and Seung, H.S. (2014). A genetic and computational approach to structurally classify neuronal types. *Nat Commun* 5, 3512.
- Sunkin, S.M., Ng, L., Lau, C., Dolbeare, T., Gilbert, T.L., Thompson, C.L., Hawrylycz, M., and Dang, C. (2013). Allen Brain Atlas: an integrated spatio-temporal portal for exploring the central nervous system. *Nucleic acids research* 41, D996–D1008.
- Susaki, E.A., Tainaka, K., Perrin, D., Kishino, F., Tawara, T., Watanabe, T.M., Yokoyama, C., Onoe, H., Eguchi, M., Yamaguchi, S., et al. (2014). Whole-brain imaging with single-cell resolution using chemical cocktails and computational analysis. *Cell* 157, 726–739.
- Suzek, B.E., Huang, H., McGarvey, P., Mazumder, R., and Wu, C.H. (2007). UniRef: comprehensive and non-redundant UniProt reference clusters. *Bioinformatics* 23, 1282–1288.
- Tanaka, N.K., Endo, K., and Ito, K. (2012). Organization of antennal lobe-associated neurons in adult *Drosophila melanogaster* brain. *Journal of Comparative Neurology* 520, 4067–4130.
- Tanaka, N.K., Tanimoto, H., and Ito, K. (2008). Neuronal assemblies of the *Drosophila* mushroom body. *Journal of Comparative Neurology* 508, 711–755.
- von Philipsborn, A.C., Liu, T., Yu, J.Y., Masser, C., Bidaye, S.S., and Dickson, B.J. (2011). Neuronal control of *Drosophila* courtship song. *Neuron* 69, 509–22.
- Ward Jr, J.H. (1963). Hierarchical grouping to optimize an objective function. *Journal of the American statistical association* 58, 236–244.
- Wong, A.M., Wang, J.W., and Axel, R. (2002). Spatial representation of the glomerular map in the *Drosophila* protocerebrum. *Cell* 109, 229–41.
- Wu, C.H., Apweiler, R., Bairoch, A., Natale, D.A., Barker, W.C., Boeckmann, B., Ferro, S., Gasteiger, E., Huang, H., Lopez, R., et al. (2006). The Universal Protein Resource (UniProt): an expanding universe of protein information. *Nucleic acids research* 34, D187–D191.
- Yorozu, S., Wong, A., Fischer, B.J., Dankert, H., Kernan, M.J., Kamikouchi, A., Ito, K., and Anderson, D.J. (2009). Distinct sensory representations of wind and near-field sound in the *Drosophila* brain. *Nature* 458, 201–205.
- Yu, H.H., Kao, C.F., He, Y., Ding, P., Kao, J.C., and Lee, T. (2010a). A complete developmental sequence of a *Drosophila* neuronal lineage as revealed by twin-spot MARCM. *PLoS Biol* 8.

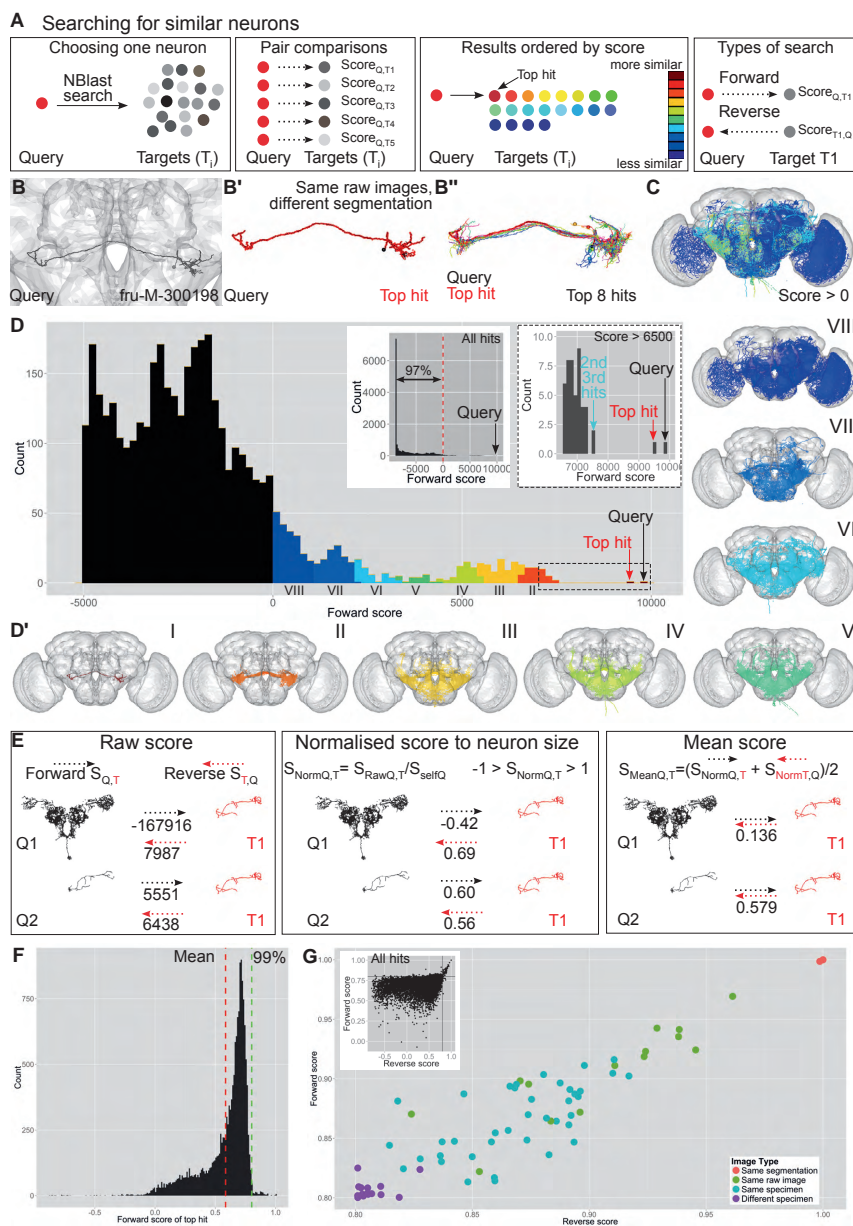
- Yu, J.Y., Kanai, M.I., Demir, E., Jefferis, G.S.X.E., and Dickson, B.J. (2010b). Cellular organization of the neural circuit that drives *Drosophila* courtship behavior. *Curr Biol* *20*, 1602–14.
- Zhu, S., Chiang, A.S., and Lee, T. (2003). Development of the *Drosophila* mushroom bodies: elaboration, remodeling and spatial organization of dendrites in the calyx. *Development* *130*, 2603–2610.
- Zingg, B., Hintiryan, H., Gou, L., Song, M.Y., Bay, M., Bienkowski, M.S., Foster, N.N., Yamashita, S., Bowman, I., Toga, A.W., et al. (2014). Neural networks of the mouse neocortex. *Cell* *156*, 1096–1111.





**Figure 1: Image preprocessing, registration and similarity score (NBLAST) algorithm**

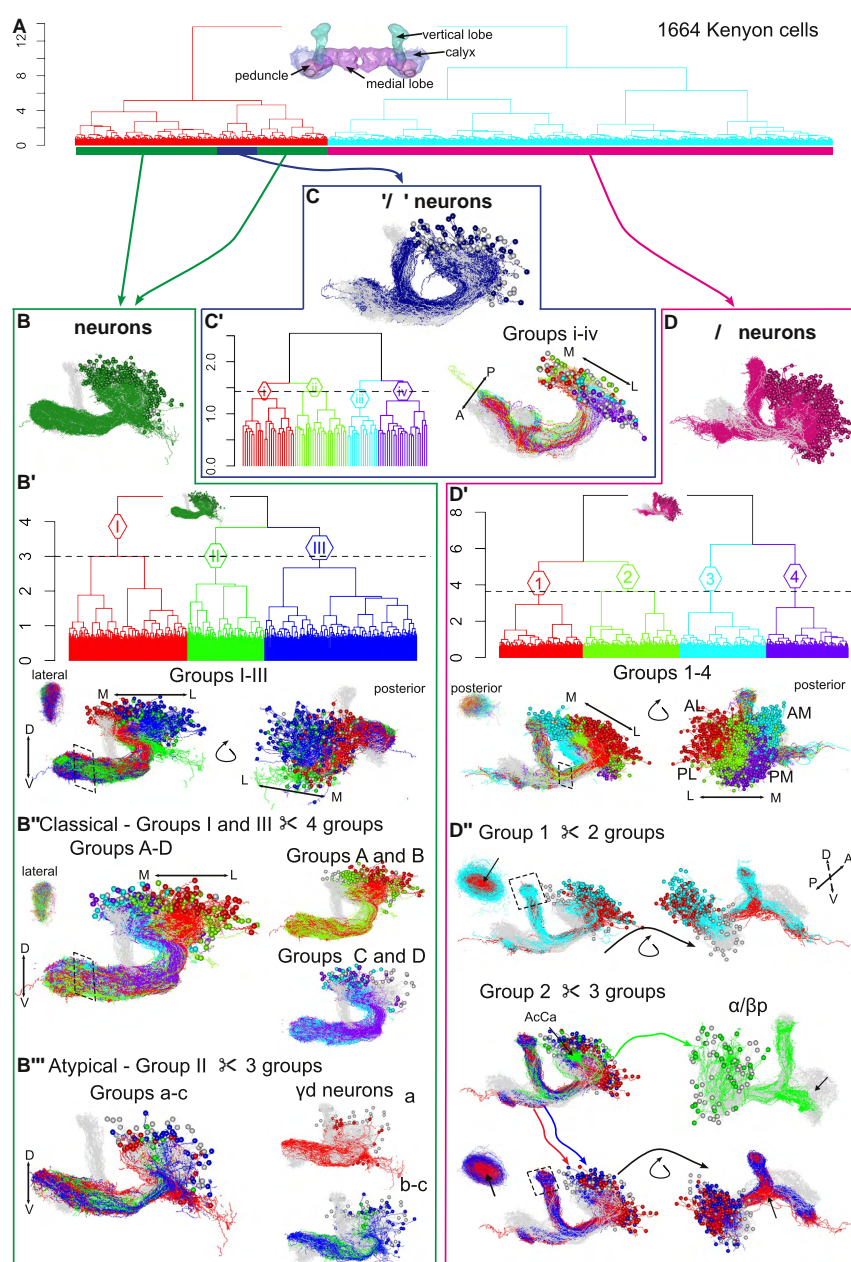
**(A)** Flowchart describing the image preprocessing and registration procedure. FlyCircuit images were split into 3 channels. The Dlg-stained brain (discs large 1) (channel 1) images were registered against the FCWB template. Registration success was assessed by comparing the template brain with the brain images after applying the registration (reformatted). The neuron image in channel 3 was skeletonized and reformatted onto the template brain. The neuron skeleton was converted into points and vectors. **(B)** After image registration, neurons on the right side of the brain were flipped onto the left side. For most neurons this was done automatically. On the left, brain plot showing 50 random neurons before and after flipping. On the right, cases for which the neuron flipping was assessed manually. These included cases in which the cell body was on or very close to the midline, with or without small primary ipsilateral neurites. **(C)** NBLAST algorithm. The similarity of two neurons (query and target), is given by a function of the distance and absolute dot product between the nearest neighbor points of the query/target pair. This distance function reflects the probability of a match between a pair of points ( $p_{match}$ ), relative to any two random points ( $p_{rand}$ ). **(D)** Diagram illustrating how nearest neighbor points are calculated. For a query (N1)/target (N2) pair, each point of N1 ( $u_i$ ) is paired to the N2 point ( $v_i$ ) that minimizes the distance ( $d_i$ ) between the points. **(E)** Calculating the distance function. Two groups of neurons were used to calculate the distribution probabilities of matching and non-matching pairs. The first corresponds to a known class of uniglomerular olfactory projection neurons (uPNs), DL2 uPNs that had been previously identified in the dataset. The second group corresponds to all remaining neurons. Random pairs of neurons were compared within each group. **(F)** Brain plot showing all DL2 neurons in the dataset that were used for this analysis. **(G)** Calculation of the distribution for matching and non-matching pairs of segments. For all segment pairs of all neuron pairs of each group, the distance and absolute dot product were plotted in a distance histogram. The distribution probability for matching ( $p_{match}$ ) or non-matching pairs ( $p_{rand}$ ) was calculated by normalizing the distance histogram to 1. When calculating the distance function,  $1 \times 10^6$  was added to both  $p_{match}$  and  $p_{rand}$  to avoid a 0 denominator. **(H)** Plot showing that the similarity score depends on the spatial location of the points (distance between points) and the direction of the vectors (absolute dot product). The score is the highest for a distance of 0  $\mu m$  and an absolute dot product of 1.



**Figure 2: Example of searching for similar neurons using NBLAST**

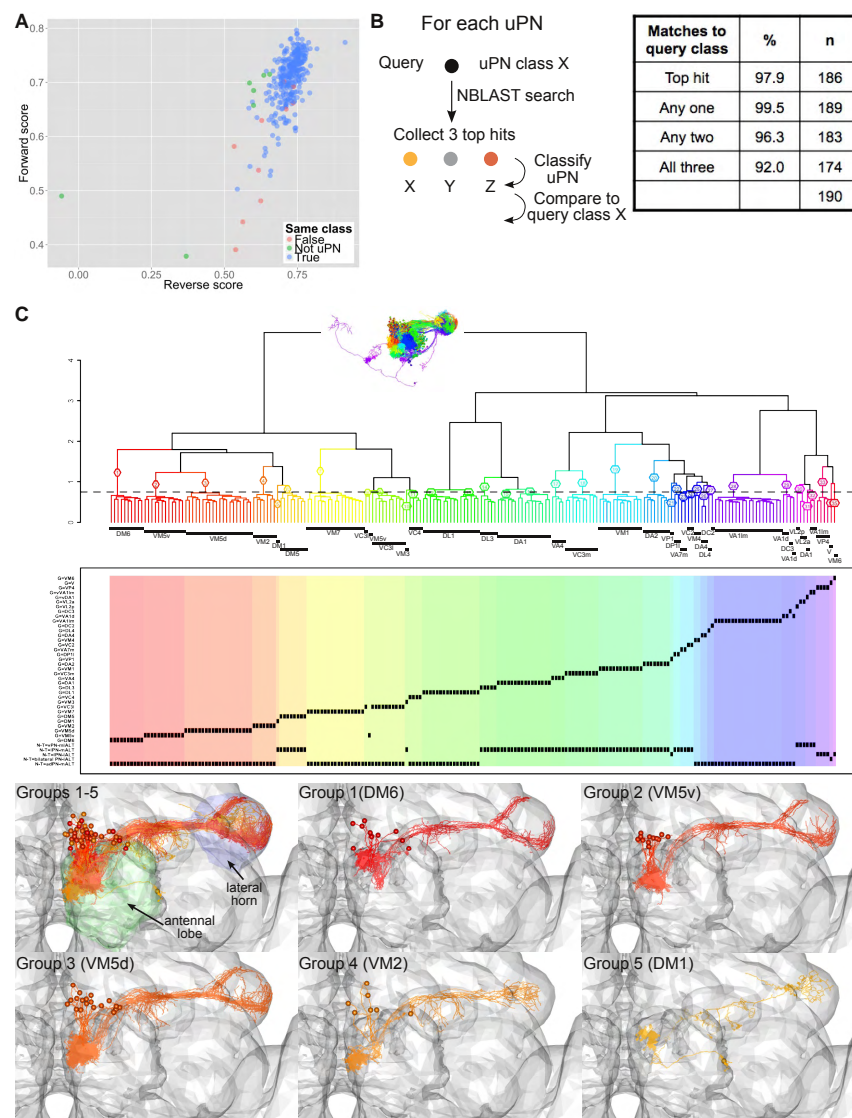
**(A)** Searching for similar neurons with NBLAST. Pair comparisons between the chosen query neuron and remaining neurons in the dataset return a similarity score, allowing the results to be ordered by similarity. A forward search returns the similarity score for the query compared to the target, whereas a reverse search returns the score of the target against the query. **(B)** NBLAST search with fru-M-300198 as query (black). Neuron plot of the query neuron. **(B')** Neuron plot of the query (black) and top hit (red). The top hit corresponds to a different segmentation of the query neuron, from the same raw image. The differences between these two images is due to minor differences during the segmentation. **(B'')** Neuron plot showing the top 8 hits. There are differences in neurite branching, length and position. **(C)** All hits with a forward score over 0, colored by score, as shown in D. **(D)** Histogram of scores for a forward search with fru-M-300198 as query. Only hits with scores over -5,000 are shown. The left inset shows the histogram of scores for all search hits. The right inset shows a zoomed view of the top hits (score > 6,500). For more examples see S1. **(D')** Neuron plots corresponding to the score bins in D. **(E)** Comparison of the raw, normalized and mean score, for two pairs of neurons: one of unequal (Q1, T1) and one of similar size (Q2, T1). The value of the raw score depends on the size of the neuron, whereas the normalized score corrects for it, by dividing the raw score by the query self-score (maximum score). Normalized scores are between -1 and 1. Mean scores are the average between the normalized forward and reverse score for a pair of neurons. These scores can be compared for different searches. **(F)** Histogram of the score for the top hit for each neuron in the whole dataset. The mean and 99th percentile are shown as a dashed red and green lines, respectively. **(G)** Plot of reverse and forward normalized scores for 72 pairs of neurons for which both the forward and reverse scores are higher than 0.8. These pairs were classified into four categories, according to the relationship between the two images: images correspond to a segmented image that is duplicated ('Same segmentation'); images correspond to different neuron segmentations from the same raw image ('Same raw image'); images correspond to two different segmented images from the same brain ('Same specimen'); images correspond to segmented images of the same or similar neurons in different brains ('Different specimen'). The inset plot shows the normalized reverse and forward scores for all top hits. The threshold of 0.8 is indicated by two black lines.





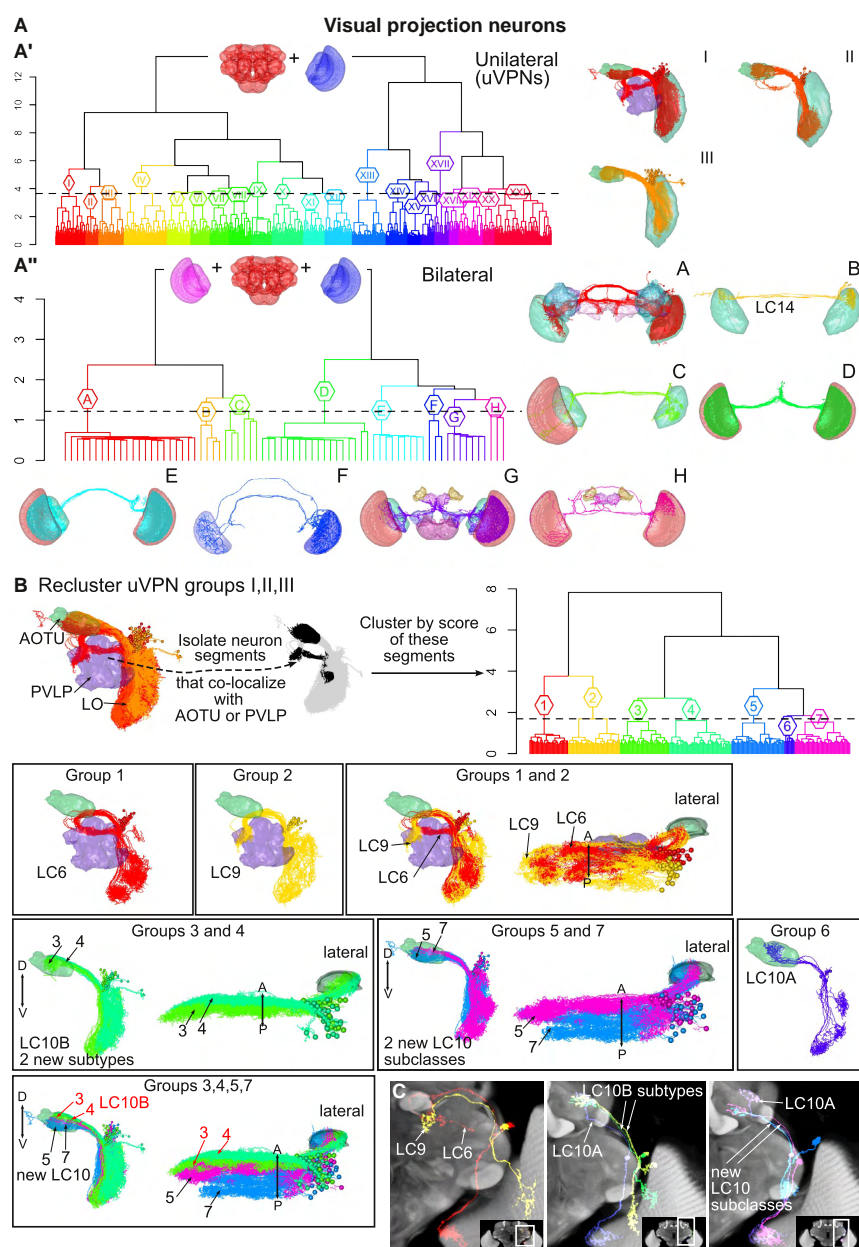
**Figure 3: Hierarchical clustering reveals Kenyon cell subtypes**

**(A)** Hierarchical clustering (HC) of Kenyon cells ( $n=1664$ ), divided into two groups. Bars below the dendrogram indicate the neurons corresponding to a specific neuron type:  $\gamma$  (in green),  $\alpha'/\beta'$  (in blue) and  $\alpha/\beta$  neurons (in magenta),  $h=8.9$ . Inset shows the mushroom body neuropil. **(B)** Neuron plot of the  $\gamma$  neurons. **(B')** HC of the  $\gamma$  neurons divided into three groups (I-III),  $h=3$ . Inset on the dendrogram shows the  $\gamma$  neurons (same as in B). Neuron plots of groups I to III. A lateral oblique and a posterior view of the neurons are shown. There are differences between the 4 groups in the calyx in the medial/lateral axis and in the dorsal/ventral axis in the  $\gamma$  lobe: the more medial group 1 is the most dorsal in the  $\gamma$  lobe. **(B'')** HC of the classic  $\gamma$  neurons, corresponding to groups I and III in B', divided into four groups (A-D). Neuron plots of groups A-D, A-B and C-D. There are differences between the 4 groups in the calyx in the medial/lateral axis and in the dorsal/ventral axis in the  $\gamma$  lobe. **(B''')** HC of the atypical  $\gamma$  neurons corresponding to group II in B', divided into three groups (a-c). Neuron plots of groups a-c, a, and b-c. Group a corresponds to  $\gamma_d$  neurons which innervate the dorsal most region of the gamma lobe and extend dendrites laterally. **(C)** Neuron plot of the  $\alpha'/\beta'$  neurons. **(C')** HC of the  $\alpha'/\beta'$  neurons, divided into four groups (i-iv),  $h=1.43$ . The groups i and iv take a more anterior route in the peduncle and  $\beta'$  lobe than groups ii and iii. Dorsolateral view is shown. **(D)** Neuron plot of the  $\alpha/\beta$  neurons. **(D')** HC of the  $\alpha/\beta$  neurons, divided into four groups (1-4),  $h=3.64$ . Inset on the dendrogram shows the  $\alpha/\beta$  neurons (same as in D). Neuron plots of groups A to D. Lateral oblique, posterior view and posterior view of a peduncle slice of these groups are shown. There are differences between the 4 groups in the calyx and in the medial/lateral axis, with each group corresponding to the indicated neuroblast clone (AM, AL, PM, PL). **(D'')** HC of groups 1 and 2. Lateral oblique, posterior oblique and a dorsal view of a peduncle slice views are shown. HC of group 1 divided into 2 subgroups. This separated the neurons into peripheral (cyan) and core (red) in the  $\alpha$  lobe. Peripheral neurons occupied a more lateral calyx position and were dorsal to core neurons in the peduncle and  $\beta$  lobe. Similar analysis to groups 3 and 4 is shown in Figure S3A. HC of group 2 divided into 3 subgroups. The red and blue subgroups match the core and peripheral neurons, respectively; the green subgroup the  $\alpha/\beta$  posterior subtype ( $\alpha/\beta_p$ ). These neurons innervate the accessory calyx and their axons terminate before reaching the most medial region of the  $\beta$  lobe. AcCa: accessory calyx. Neurons in grey: Kenyon cell exemplars.



**Figure 4: Hierarchical clustering of uniglomerular olfactory projection neurons**

**(A)** Plot of the reverse and forward normalized score for the top hit in a NBLAST search using the uniglomerular olfactory projection neurons (uPNs) as queries. Only uPN types for which we have more than one example and unique query/target pairs are included in this analysis (n=327). For each query neuron, we identified the cases for which both the top hit and query were of the same class (True) (n=319); the top hit is a uPN but does not match the class of the query (False) (n=4), or the top hit is not a uPN (Not uPN) (n=4). **(B)** The top three hits for each uPN (mean score) were collected and the neuron type of the hits and query was compared. Only non-DL2 uPNs for which we had more than three neurons examples were used (n=190). **(C)** Hierarchical clustering of uPNs (non-DL2s) (n=214) divided into 35 groups (1–35),  $h=0.75$ . Dendrogram and picket plot, showing the glomerulus, tract/lineage for each neuron. The neuron plot inset shows the uPNs colored by dendrogram group. Below the leaves, the number of neurons that innervate each glomerulus is indicated by the black rectangles. On the left of the picket plot, the neuron types characterized by glomerulus (G), lineage (N) or tract (T) are shown. Neurons that innervate DA1 and VA1Im glomeruli but originate from the ventral lineage instead of the lateral or anterodorsal, respectively, are indicated as vVA1Im and vDA1 on the picket plot labels. The dendrogram groups correspond to single and unique neuron types except for DL1 and DA1 neurons which are split into 2 groups. Below the picket plot, neuron plots corresponding to dendrogram groups 1–5 and to each of these individual groups, colored by dendrogram group. In the first, the antennal lobe is in green, the lateral horn in purple. vPN: ventral lineage; adPN: anterodorsal lineage; lPN: lateral lineage; mALT: medial antennal lobe tract; mlALT: mediolateral antennal lobe tract; lALT: lateral antennal lobe tract.

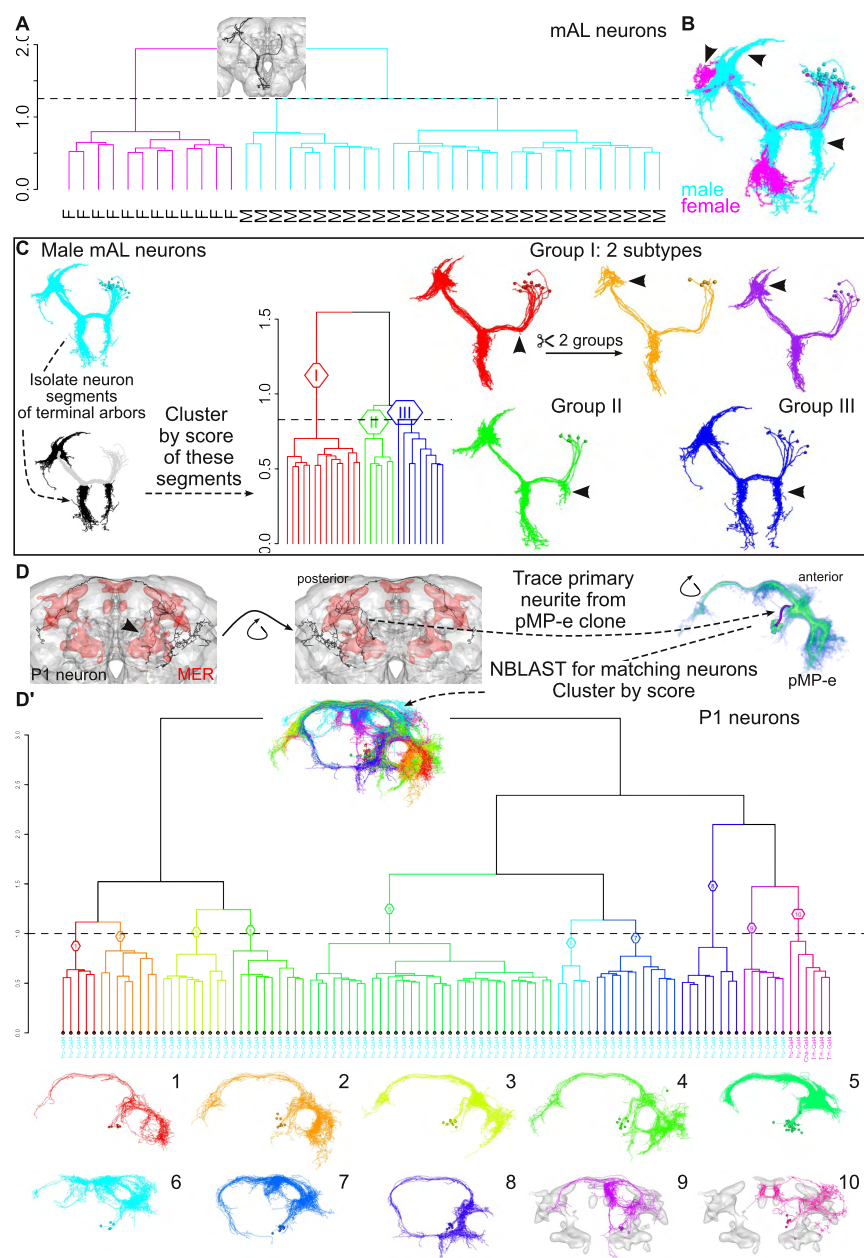


**Figure 5: Similarity search and hierarchical clustering of visual projection neurons**

**(A)** Clustering analysis of unilateral (uVPNs) and bilateral visual projection neurons (bVPNs), defined as neurons with segments that overlap one or two optic lobes, respectively, and some central brain neuropil. **(A')** Hierarchical clustering (HC) of unilateral visual projection neurons (uVPNs), divided into 21 groups (I–XXI),  $h=3.65$ . Inset on the dendrogram shows the neuropils considered for the overlap. To the right, neuron plots of groups I to III. The neuropils that contain the most overlap are shown. Other neuron plots are shown in Figure S4. **(A'')** HC of bVPNs, divided into 8 groups (A–H),  $h=1.22$ . Inset on the dendrogram shows the neuropils considered for the overlap. To the right and below, neuron plots of dendrogram groups. Group 2 corresponds to the LC14 neuron type that connects the 2 lobulas, with one outlier terminating in the medulla. The neuropils that contain the most overlap are shown.

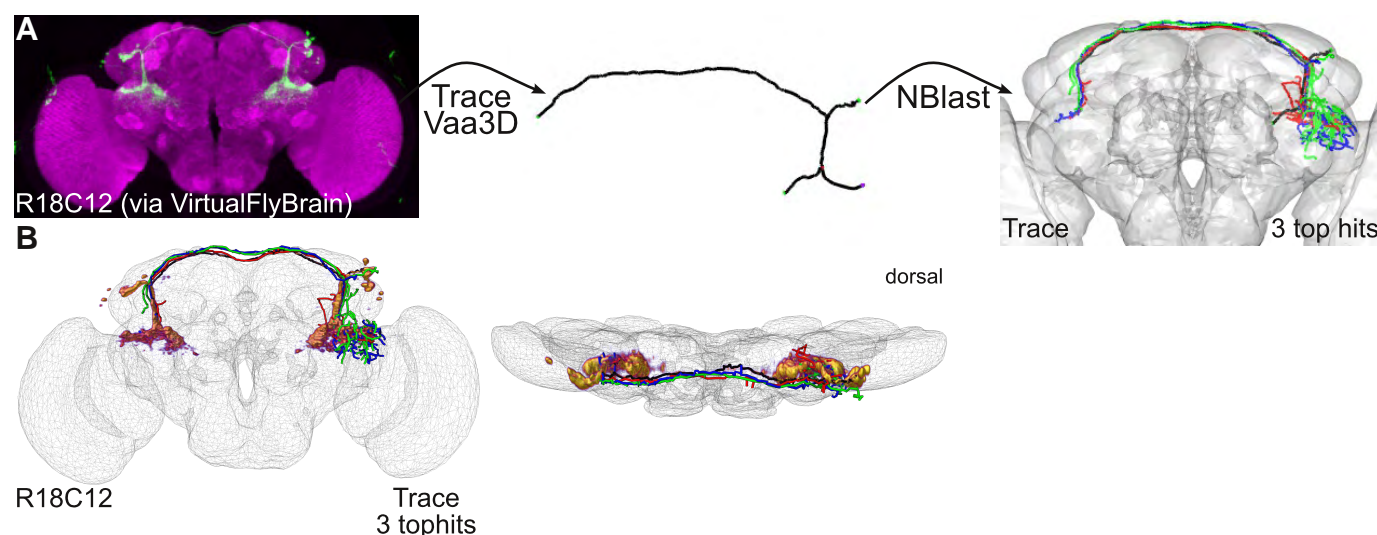
**(B)** Reclustering of uVPN groups I, II and III from A'. These neurons have dorsal cell bodies, arborize in the lobula (LO) and project to anterior optic tubercle (AOTU) via the anterior optic tract or to the posterior ventrolateral protocerebrum (PVLP). The neuron segments that co-localize with either the AOTU or PVLP were isolated, followed by HC of the neurons based on the NBLAST score of these neuron segments. The dendrogram was divided into seven groups (1–7),  $h=1.69$ . Neuron plots corresponding to the dendrogram groups. An anterior and a lateral view are shown. Some of dendrogram groups were matched to known uVPN types. Group 1 corresponds to LC6 neurons, group 2 to LC9. These 2 groups innervate the PVLP, and show some differences in the lobula lamination in the anterior/posterior axis. Groups 3 and 4 seem to correspond to two new subtypes of LC10B, that innervate the dorsal AOTU. They show a clear distinction in AOTU and lobula lamination, with group 4 being the dorsalmost in the AOTU and the most anterior in the lobula. Groups 5 and 7 are possible new subclasses of LC10, that innervate the ventral AOTU. They show a clear distinction in AOTU and lobula lamination, with group 7 being the dorsalmost in the AOTU (but ventral to group 3) and the most anterior in the lobula (but posterior to group 3). Group 6 corresponds to LC10A neurons that project through the ventral AOTU and turn sharply dorsally in the middle region. **(C)** Overlay of Z projections of registered image stacks of example neurons from the types identified in B on a partial Z projection of the template brain (a different one for each panel). The white rectangle on the inset shows the location of the zoomed in area. LC: lobula columnar neuron.





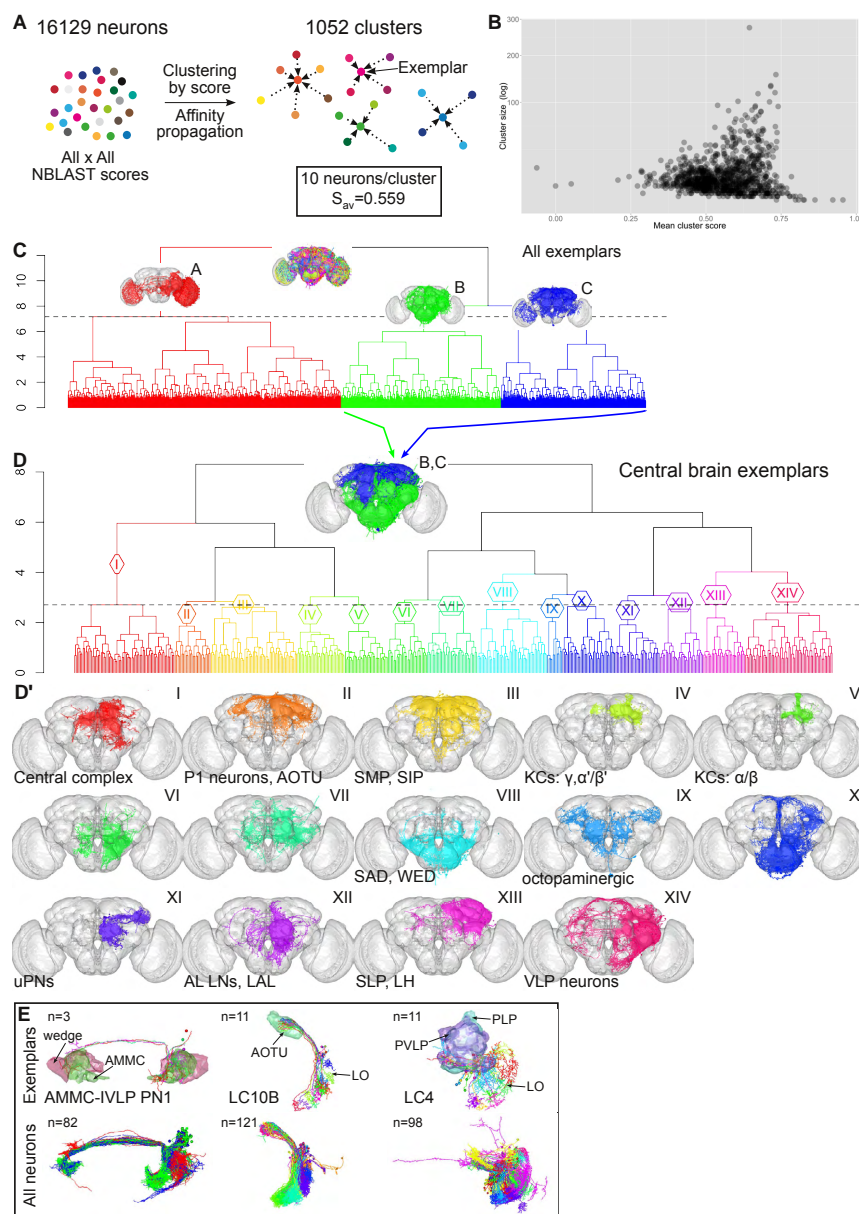
**Figure 6: Hierarchical clustering of *fruitless*-expressing mAL and P1 neurons**

**(A)** Analysis of the mAL neurons. Hierarchical clustering (HC) of the hits, divided into 2 groups ( $h=1.25$ ). The mAL neuron used as the Nblast query, fru-M-500159, is shown in the inset. Hits with a normalized score over 0.2 were collected. The leaf labels indicate the gender of the neuron: 'F' for female and 'M' for male. **(B)** Neuron plot of the 2 dendrogram groups corresponding to male (in cyan) and female (in magenta) mAL neurons. **(C)** Analysis of the male mAL neurons. The neuron segments corresponding to the terminal arbors (ipsi- and contralateral) were isolated and the neurons were clustered based on the score of these segments. HC of neurons, divided into 3 groups (groups I–III) ( $h=0.83$ ), that reflect differences in the length of the ventral ipsilateral branch (arrowhead). Group I can be further subdivided into two different subtypes, which differ in the shape and extent of their dorsal contralateral arborisation (arrowhead). **(D)** Analysis of the P1 neurons. Neuron plot of a P1 neuron, fru-M-400046. The male enlarged region (MER) is shown in red. Anterior and posterior views are shown. Volume rendering of the pMP-e *fruitless* neuroblast clone, which gives rise to P1 neurons. The distinctive primary neurite was traced and used on a NBLAST search for matching neurons. **(D')** HC of hits for a search against the P1 primary neurite divided into 10 groups (1–10) ( $h=1$ , indicated by dashed line). This group of neurons corresponds to a subset of neurons obtained after a first HC analysis. Hits with a normalized score over 0.25 were collected and further selected. The inset shows a neuron plot with groups 1–10. The leaf labels show the GAL4 driver used to obtain that neuron; the colors follow the gender: cyan for male and magenta for female. Below the dendrogram, neuron plots of each group. The MER is shown in grey for groups 9 and 10.



**Figure 7: Identifying neuron types from GAL4 traces**

**(A)** Maximum Z projection of FlyLight line R18C12 registered to JFRC2. The image was downloaded from the Virtual Fly Brain website. The neuron on the left was traced using Vaa3D, its coordinates adjusted to the correct JFRC2 voxel dimensions and transformed into FCWB space. Neuron plot of the top 3 hits of an NBLAST search with normalised scores against all FlyCircuit neurons. The 3 top hits are very similar to the traced neuron. **(B)** Volume rendering of the traced neuron, R18C12 and 3 top hits in JFRC2 space. Neuron plots from a frontal and dorsal views are shown. Trace in black; 1st top hit in red; 2nd top hit in green; 3rd top hit in blue.



**Figure 8: Affinity propagation clustering**

**(A)** Clustering by affinity propagation. This method uses the all-by-all matrix of NBLAST scores for the 16,129 neurons. This method defined exemplars, which are representative members of each cluster. An affinity propagation clustering of the dataset generated 1,052 clusters, with an average of 10 neurons per cluster and a similarity score of 0.559. **(B)** Plot showing the mean cluster score versus cluster size. **(C)** Hierarchical clustering (HC) of the 1,052 exemplars, dividing them into three groups (A–C). Group A corresponds mostly to optic lobe and VPN neurons; groups B and C to central brain neurons. The insets on the dendrogram show the neuron of these groups. The main neuron types or innervated neuropils are noted. **(D)** HC of central brain exemplars (groups B and C, inset on dendrogram), divided into 14 groups,  $h=2.7$ . **(D')** Neurons corresponding to the dendrogram groups in D. **(E)** Affinity propagation clusters of defined neuron types. Neuron plot of exemplars (top row) or all neurons (bottom row) for auditory AMMC-IVLP PN1 neurons (compare with Figure S5D) and VPN types LC10B (compare with Figure 5B) and LC4 (compare with Figure S4B). The number of exemplars and neurons is indicated on the top left corner for each example. The AMMC is shown in green, the wedge in magenta. AMMC: antennal mechanosensory and motor center; AOTU: anterior optic tubercle; LO: lobula; PVLP: posterior ventrolateral protocerebrum; PLP: posterior lateral protocerebrum.



# Supplemental Information

## Supplemental Results

### Kenyon cell analysis

In order to understand if the relative position of the classical  $\gamma$  neurites was maintained in between the calyx and the  $\gamma$  lobe, we clustered the neurons based on the scores of the segments in the peduncle. We took neuron skeletons from classical  $\gamma$  neurons and isolated the axon arbors that co-localized with the peduncle volume (Figure S3B'). We then carried out a new clustering based on all-by-all NBLAST scores of these partial skeletons, cutting the dendrogram at a level defined by visual inspection (4 groups). The overall organization almost fully recapitulated the positioning of the neurites in the whole neuron analysis (compare Figure 3B'–B'' with Figure S3B'). A clear and expected lamination was found in the peduncle, with neurites occupying the most outer stratum. Differences in the medial to lateral positioning of neurites in the calyx followed the previously observed organization, with the most medial groups occupying the dorsal region of the gamma lobe.

In order to investigate the stereotypical organization of  $\alpha/\beta$  neurites, we performed a similar analysis as for the classic  $\gamma$  neurons, isolating the axon arbors that co-localized with the peduncle for groups 1 to 4 (Figure S3B''). The new clustering based on peduncle position of these partial neuron skeletons did not recapitulate the relative positions of the calyx neurites for each of the neuroblast clones observed in the whole neuron analysis (compare Figure 3D' with Figure S3B''). In addition, there was no clear organization of neurites in the  $\alpha$  lobe that correlated with their position in the peduncle.

## Supplemental Figures and Table

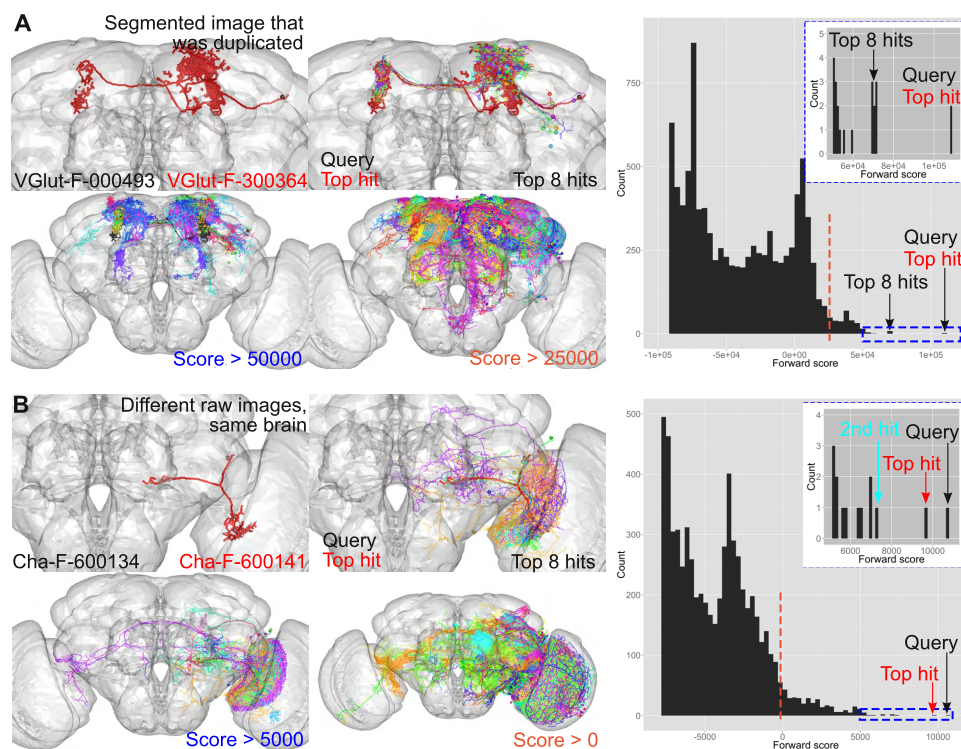


Figure S1: Neuron search with NBLAST

(A) NBLAST search with VGlut-F-000493 as query (black). Neuron plots of (from higher to lower score): the query and top hit, top 8 hits, hits with a score over 50,000 and hits with a score over 25,000. The top hit corresponds to a segmented image that was duplicated. It perfectly overlays the query neuron. As the score decreases, so does the similarity of the hits to the query. On the right, histogram of forward scores. Only hits with scores over  $-100,000$  are shown. The score of the query, top hit and top 8 hits are indicated. A dashed orange line marks 0. The left inset shows a zoomed view of the top hits (score > 50,000) (dashed blue rectangle in main plot). The score of the query, top hit and top 8 hits are indicated. (B) NBLAST search with Cha-F-600134 as query (black). Neuron plots of (from higher to lower score): the query and top hit, top 8 hits, hits with a score over 5,000 and hits with a score over 0. The top hit corresponds to an image of a neuron from the same brain but from a different raw image. It is very similar to the query neuron. As the score decreases, so does the similarity of the hits to the query. On the right, histogram of forward scores. Only hits with scores over  $-8,000$  are shown. The score of the query and top hit are indicated. A dashed orange line marks 0. The left inset shows a zoomed view of the top hits (score > 5,000) (dashed blue rectangle in main plot). The score of the query, top hit and second top hits are indicated.

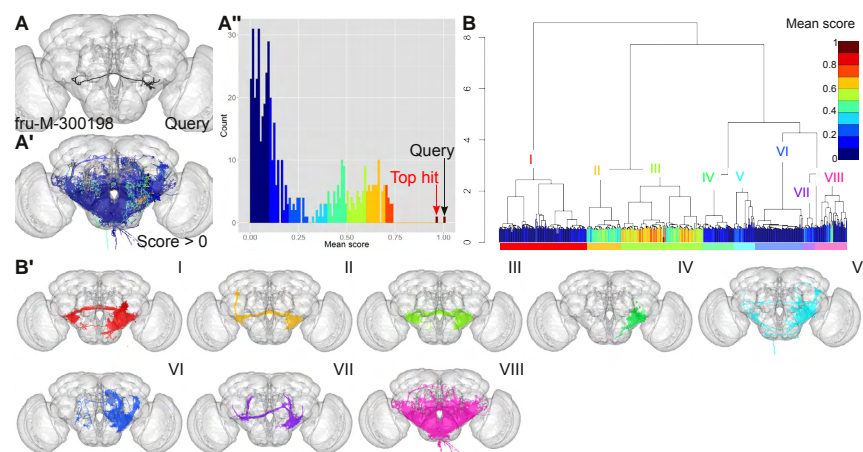
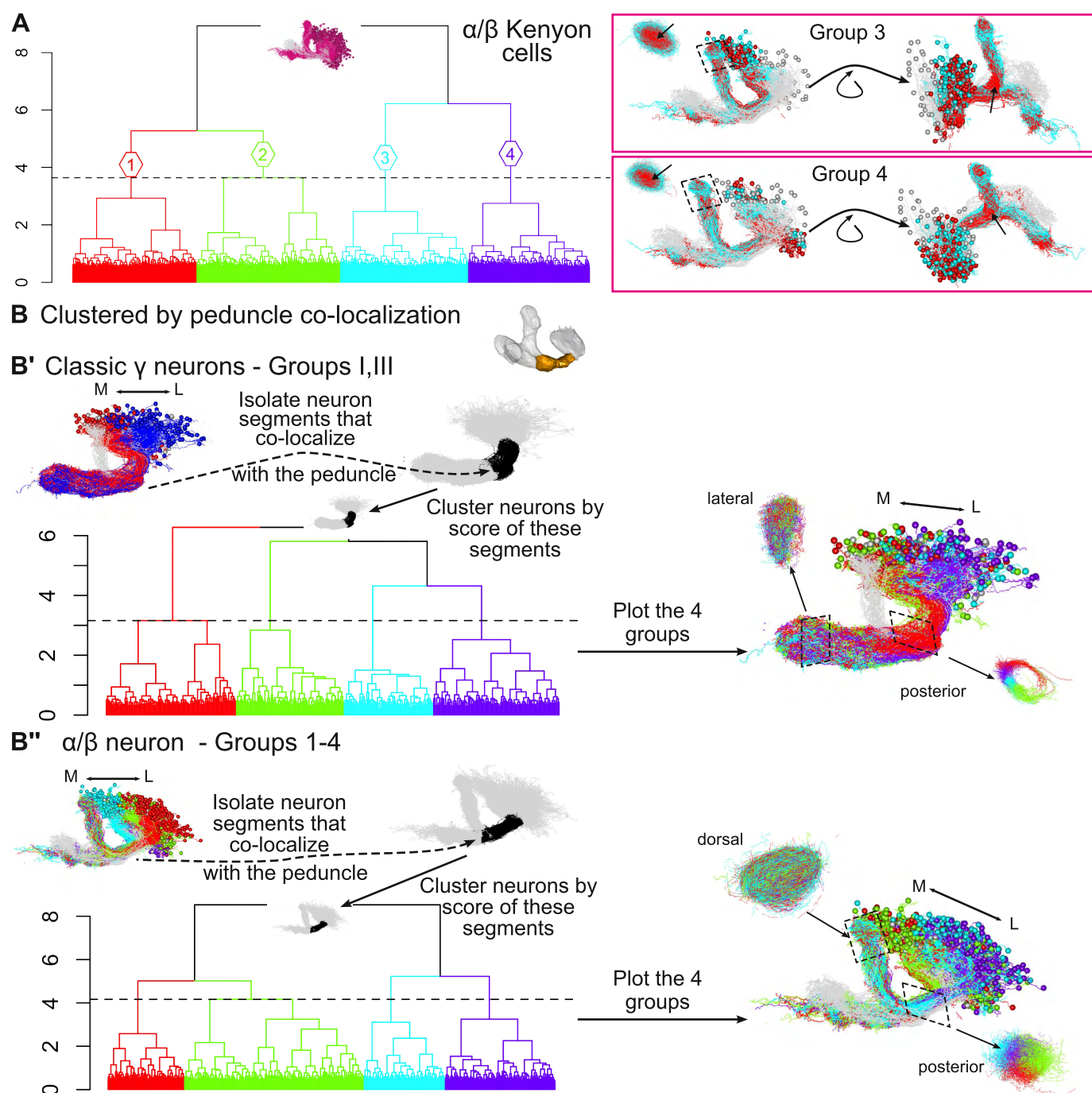


Figure S2: **NBLAST search using mean scores**

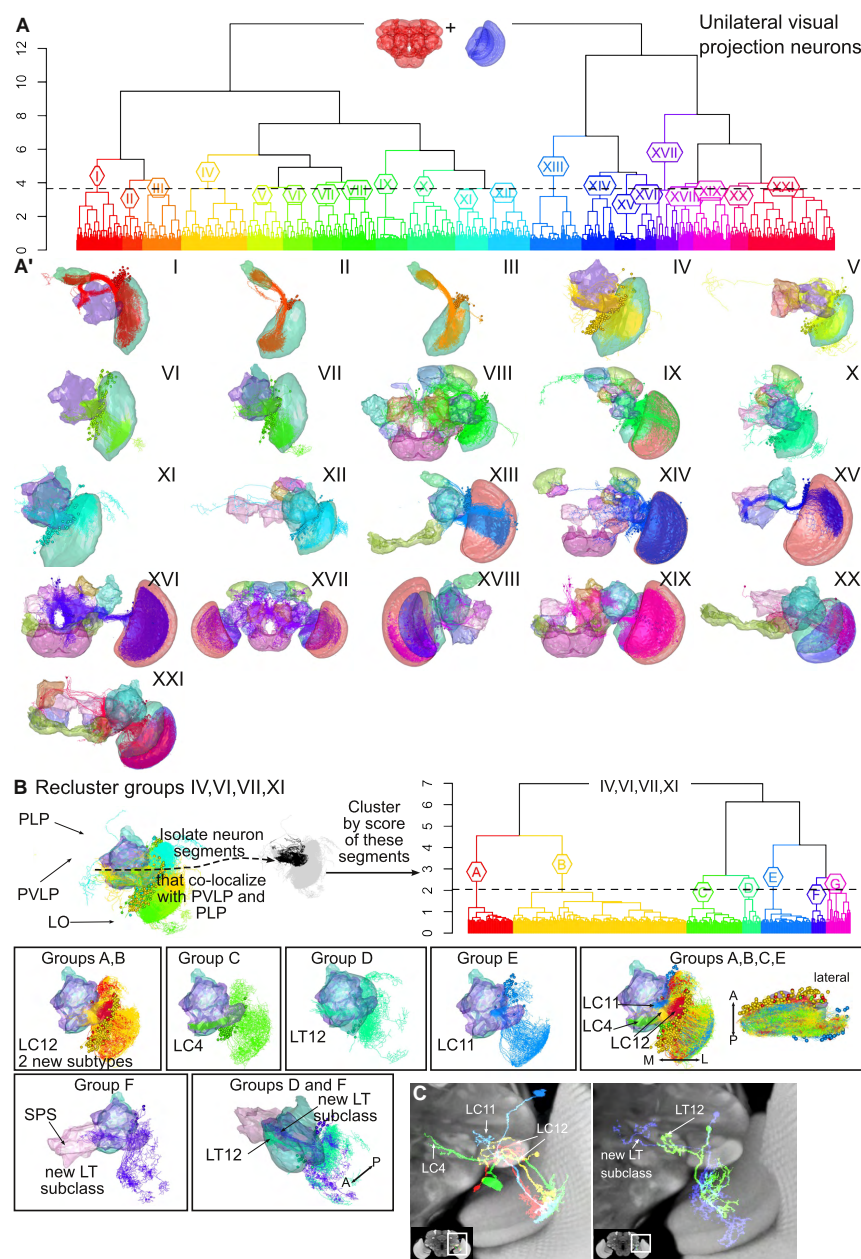
(A) The query neuron fru-M-300198 (same as inFigure 2B). (A') Neuron plot of the hits with a mean score over 0 for a search against the query. Hits are colored by score bin (10), as in A''. (A'') Histogram of mean scores for hits against fru-M-300198 with a score over 0 divided into 10 bins (indicated in the scale bar in B). (B) Hierarchical clustering of hits with a mean score over 0. The leaves of the dendrogram are colored by score (same as in A''), and as shown in the scale bar. The dendrogram was divided into eight groups (I-VIII), with each one being assigned a color, shown on the colored rectangle below the leaves. The query neuron is in group III, and the hits with the higher scores are in groups II and III. (B') Neuron plots corresponding to the dendrogram groups (I-VIII), following the colors assigned to each group. Groups II and III, corresponding to the highest scores, are the most similar neurons to the query.



**Figure S3: Clustering analysis of Kenyon cells**

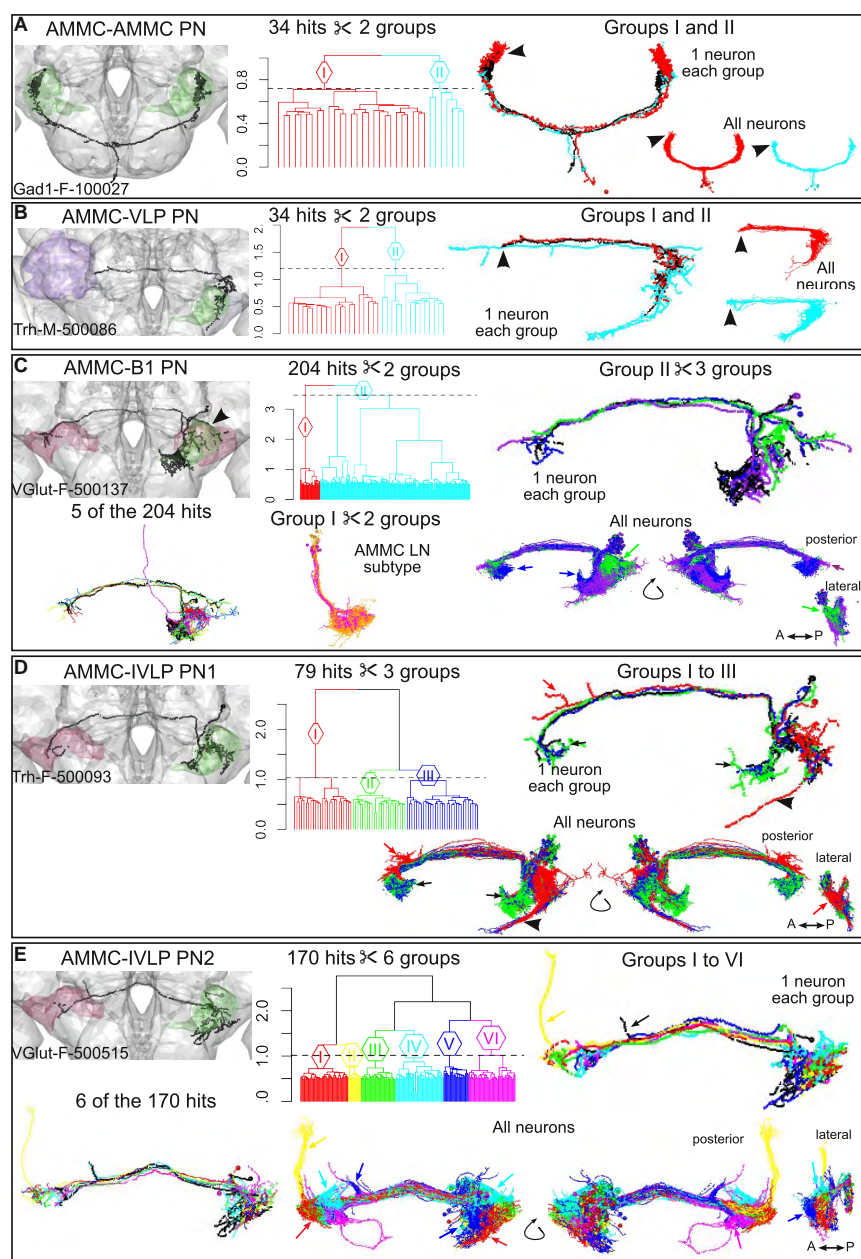
**(A)** Hierarchical clustering (HC) of the  $\alpha/\beta$  neurons, divided into four groups (1–4) ( $h=3.64$ ) (same as in Figure 3D'). The inset on the dendrogram shows the  $\alpha/\beta$  neurons. Groups 3 and 4 were clustered and divided into 2 groups each. This separates the neurons into peripheral (cyan) and core (red) in the  $\alpha$  lobe. Peripheral neurons occupied a more lateral calyx position and were dorsal to core neurons in the peduncle and  $\beta$  lobe. Similar analysis to group 1 is shown in Figure 3D'. Lateral oblique, posterior oblique and a dorsal view of a peduncle slice (position indicated by dashed rectangle) are shown. **(B)** Reclustering of Kenyon cells based on the co-localization of neurons segments in the peduncle. The neuron segments that co-localize with the peduncle were isolated, followed by HC of the neurons based on the NBLAST score of the segments. **(B')** HC of the neurons segments of classic  $\gamma$  neurons, groups I and III (see Figure 3B') divided into 4 groups,  $h=3.16$ . Neuron plot of the 4 groups. A posterior view of a slice of the peduncle shows an expected clear organization. It correlates to the position of the neurons in the calyx, with more medial neurons (cyan and green) being dorsal and ventral in the peduncle than more lateral neurons (red and purple). No clear structural organization is discernible in a lateral view of a slice of the  $\gamma$  lobe. **(B'')** HC of the neurons segments of classic  $\alpha/\beta$  neurons, groups 1 to 4 (see Figure 3D') divided into 4 groups,  $h=4.16$ . Neuron plot of the 4 groups. A posterior view of a peduncle slice shows an expected clear organization. It correlates to the position of the neurons in the calyx, with more medial neurons (cyan and green) being ventrolateral in the peduncle than more lateral neurons (red and purple). No structural organization is discernible in a dorsal view of a slice of the  $\alpha$  lobe. For all neuron plots, the neurons in grey correspond to the Kenyon cell exemplars.





**Figure S4: Similarity search and hierarchical clustering of unilateral visual projection neurons**

**(A)** Hierarchical clustering (HC) analysis of unilateral visual projection neurons (uVPNs), defined as neurons with segments that overlap one optic lobe, and some central brain neuropil. The dendrogram was divided into 21 groups (I–XXI),  $h=3.65$ . Inset on the dendrogram shows the neuropils considered for the overlap. Below, neuron plots of groups I to XXI. The neuropils that contain the most overlap are shown. **(B)** Reclustering of uVPN groups IV, VI, VII and XI from A. These neurons arborize in the lobula (LO) and project to the posterior ventrolateral protocerebrum (PVLP) or posterior lateral protocerebrum (PLP). The neuron segments that co-localize with either the PVLP or PLP were isolated, followed by HC of the neurons based on the NBLAST score of these neuron segments. The dendrogram was divided into seven groups (A–G),  $h=2.04$ . Neuron plots corresponding to the dendrogram groups. An anterior, a lateral or lateral oblique views are shown. Some of dendrogram groups were matched to known uVPN types. Groups A and B possibly correspond to two LC12 subtypes, that innervate the more lateral PVLP glomeruli. Group B innervates a more anterior and medial glomeruli than group A (see also C). Group C corresponds to LC4, that innervates a lateral PVLP glomeruli and ventral to LC12. Group D corresponds to LT12 neurons, that project from the lateral to the medial PVLP, posterior to LC4. Group E corresponds to LC11, that innervates a lateral PVLP glomeruli, dorsal to LC12. These neurons extends along the posterior PVLP and make a sharp anterior turn, terminating with a blunt-stick like ending in the lateral PVLP. Group F corresponds to a possibly new LT subclass, with neurons projecting posteriorly to LT12 in the PVLP and extending into the superior posterior slope (SPS). **(C)** Overlay of Z projections of registered image stacks of example neurons from the types identified in B on a partial Z projection of the template brain (a different one for each panel). The white rectangle on the inset shows the location of the zoomed in area. LC: lobula columnar neuron; LT: lobula tangential neuron.



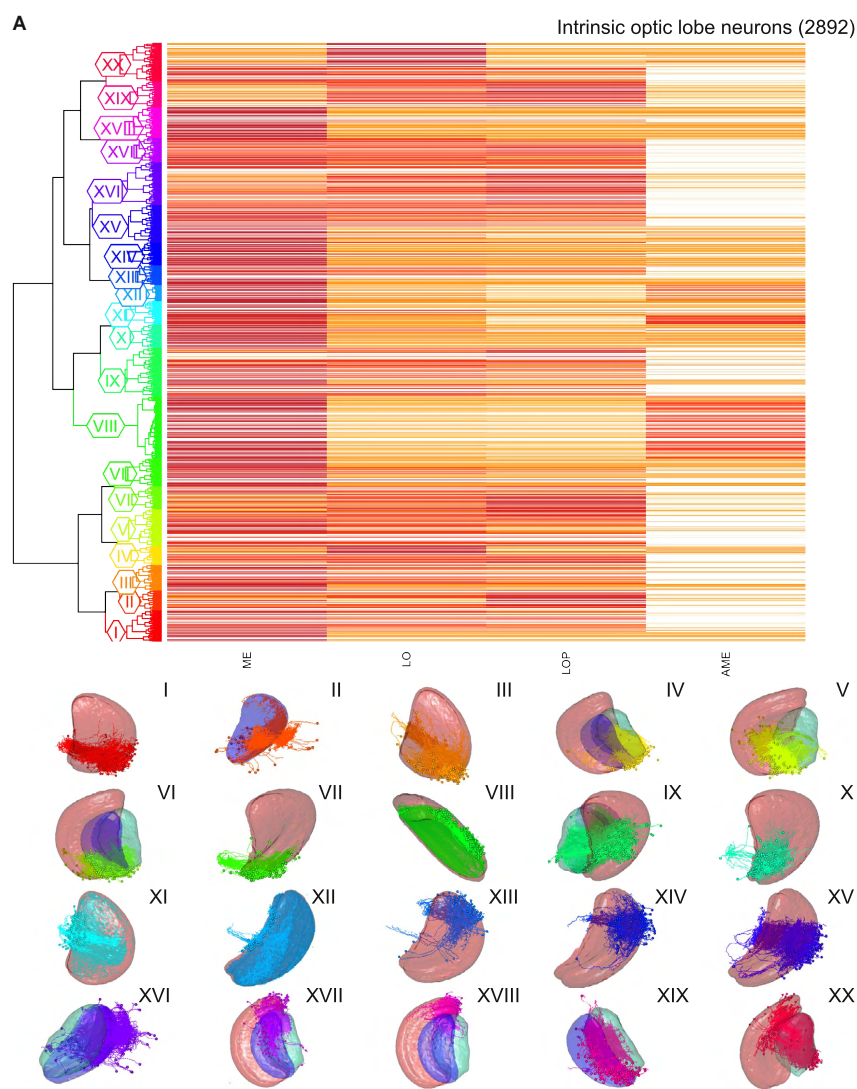
**Figure S5: Searching and classifying auditory neuron types**

Searches were done using types identified by Lai et al. (2012). A first search using the neuron named by Lai for each type (shown on the left panel with wedge in green, AMMC in magenta and PVLN in purple) was done to collect the possible candidates. A second search was then done using these neurons as queries and collecting all the high scorers (over 0.5). The dendrogram and neuron plots (anterior, posterior and lateral views) showing either one or all neurons from each clustering group are shown on the middle and right panels. **(A)** Hierarchical clustering (HC) of AMMC-AMMC projection neurons (PN). The 34 top scorers were clustered and divided into 2 groups,  $h=0.72$ . Neuron plot of the query neuron (black), and one neuron from each group (red and cyan). To the right, neuron plots of each group. Differences in anterior length of the left projection are indicated by an arrowhead. **(B)** HC of AMMC-VLP PNs. The 34 top scorers were clustered and divided into 2 groups,  $h=1.2$ . Neuron plot of the query neuron (black), and one neuron from each group (red and cyan). On the left, neuron plots of each group. **(C)** HC of AMMC-B1 PNs. Five hits of the 204 top scorers are shown on the left. The 204 top scorers were clustered and divided into 2 groups,  $h=3.46$ . Group I was matched to an unidentified type of AMMC local neurons (LN). It was clustered and divided into 2 groups,  $h=0.77$ . Group II corresponds to a mix of AMMC-B1 PNs and AMMC-IVLP PNs, the neurons were clustered and divided into 3 groups,  $h=1.5$ . Neuron plot of the query neuron (black), and one neuron from each group (purple and green). Below, neuron plots of the 3 groups. Arrows indicate differences between groups. **(D)** HC of AMMC-IVLP PN1. The 79 top scorers were clustered and divided into 3 groups,  $h=1.03$ . Neuron plot of the query neuron (black), and one neuron from each group (red, green and blue). Below, anterior, posterior and lateral view neuron plots of the 3 groups. **(E)** HC of AMMC-IVLP PN2. Six hits of the 170 top scorers are shown on the left. The 170 top scorers were clustered and divided into 6 groups,  $h=1.02$ . Neuron plot of the query neuron (black), and one neuron from each group (red, yellow, green, cyan, blue, magenta). Below, anterior, posterior and lateral view neuron plots of the 6 groups. Arrows indicate differences between groups.



Panel	Neuron type	Comments
A	AMMC-AMMC PN (2 possible subtypes)	These neurons innervate both AMMCs, with a ventral cell body. Group I extends more dorsally than group II.
B	AMMC-VLP (2 possible subtypes)	These neurons innervate the ipsilateral AMMC and the contralateral VLP. Group II extends more laterally than group I.
C	AMMC-B1 PN (3 possible subtypes)	These neurons innervate the ipsilateral AMMC and IVLP and the contralateral IVLP. Blue group innervates more medial regions and has more extensive innervation contralaterally; purple group innervates more anterior and posterior regions ipsilaterally, and the green group innervates the dorsal regions in the ipsilateral AMMC.
C	New AMMC-LN type (2 subtypes)	A type of AMMC LN, with a dorsal cell body. Two possible subtypes: the magenta group innervates more dorsal regions than the orange group.
D	AMMC-IVLP PN1 (3 possible subtypes)	These neurons innervate the ipsilateral AMMC and the contralateral IVLP. Red group has a more dorsomedial ipsilateral innervation more medial regions, with some dorsal medial branches in the contralateral hemisphere; some neurons extend a long neurite ventrally ipsilaterally. The green group innervates the more ventromedial regions ipsilaterally. The blue group is similar to the green one, although it does not extend as ventrally in the ipsilateral side, and a few neurons extend a long neurite ventrally (similar to red group).
E	AMMC-IVLP PN2 (6 possible subtypes)	These neurons innervate the ipsilateral AMMC and the contralateral IVLP, with a posterior cell body. Group I innervates the more lateral regions in both hemispheres. Group II has a long dorsal branch in the contralateral hemisphere, at the lateral edge of the neuron. Group III and IV are very similar, with the latter innervating more dorsal regions in both hemispheres. Group V corresponds to the strict definition of the neuron type by <a href="#">Lai et al. (2012)</a> , showing a short dorsal branch just medial to the contralateral IVLP. Group VI are similar to group V, with a few neurons showing a short dorsal branch, and innervating a more ventral region in the contralateral IVLP.

Table S1: Correspondence between hierarchical clustering of auditory neuron via NBLAST scores and previously determined neuron types



**Figure S6: Hierarchical clustering of intrinsic optic lobe neurons**

**(A)** Hierarchical clustering of intrinsic optic lobe neurons. This neuron set was defined as any neuron that overlapped only one of the optic lobes and with no arborization in the central brain neuropils. Dendrogram of the intrinsic optic lobe neurons, divided into into 20 groups (I–XX) with the corresponding heatmap calculated from the neuropil overlap in the different neuropils: medulla (ME), lobula (LO), lobula plate (LOP) and accessory medulla (AME). The values were log transformed. Neuron plots corresponding to the dendrogram groups are shown below. The neuropils for which the overlap is more significant are plotted. Although some organizational structure is seen, the dendrogram groups do not represent unique types.