

Manuscript Number:	GIGA-D-18-00086	
Full Title:	Clustering trees: a visualisation for evaluating clusterings at multiple resolutions	
Article Type:	Research	
Funding Information:	Department of Education, Australian Government	Mr Luke Zappia
	National Health and Medical Research Council (APP1126157)	Dr Alicia Oshlack
Abstract:	<p>Clustering techniques are widely used in the analysis of large data sets to group together samples with similar properties. For example, clustering is often used in the field of single-cell RNA-sequencing in order to identify different cell types present in a tissue sample. There are many algorithms for performing clustering and the results can vary substantially. In particular, the number of groups present in a data set is often unknown and the number of clusters identified by an algorithm can change based on the parameters used. To explore and examine the impact of varying clustering resolution we present clustering trees. This visualisation shows the relationships between clusters at multiple resolutions allowing researchers to see how samples move as the number of clusters increases. In addition, meta-information can be overlaid on the tree to inform the choice of resolution and guide in identification of clusters. We illustrate the uses of clustering trees using two examples, the classical iris dataset and a complex single-cell RNA-sequencing dataset.</p>	
Corresponding Author:	Alicia Oshlack AUSTRALIA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Luke Zappia	
First Author Secondary Information:		
Order of Authors:	Luke Zappia Alicia Oshlack	
Order of Authors Secondary Information:		
Opposed Reviewers:		
Additional Information:		
Question	Response	
Are you submitting this manuscript to a special series or article collection?	No	
Experimental design and statistics	Yes	
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available		

<p>in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

Clustering trees: a visualisation for evaluating clusterings at multiple resolutions

Luke Zappia (1, 2)

Alicia Oshlack (1, 2)

1 Bioinformatics, Murdoch Children's Research Institute; 2 School of Biosciences, University of Melbourne

Clustering techniques are widely used in the analysis of large data sets to group together samples with similar properties. For example, clustering is often used in the field of single-cell RNA-sequencing in order to identify different cell types present in a tissue sample. There are many algorithms for performing clustering and the results can vary substantially. In particular, the number of groups present in a data set is often unknown and the number of clusters identified by an algorithm can change based on the parameters used. To explore and examine the impact of varying clustering resolution we present clustering trees. This visualisation shows the relationships between clusters at multiple resolutions allowing researchers to see how samples move as the number of clusters increases. In addition, meta-information can be overlaid on the tree to inform the choice of resolution and guide in identification of clusters. We illustrate the uses of clustering trees using two examples, the classical iris dataset and a complex single-cell RNA-sequencing dataset.

Keywords: Clustering - Visualisation - scRNA-seq

Introduction

Clustering analysis is commonly used to group similar samples across a diverse range of applications. Typically, the goal of clustering is to form groups of samples that are more similar to each other than to samples in other groups. While fuzzy or soft clustering assigns each sample to every cluster with some probability, and hierarchical clustering forms a tree of samples, most methods form hard clusters where each sample is assigned to a single group. This goal can be achieved in a variety of ways, such as by considering the distances between sample (e.g. k -means¹⁻³, PAM⁴), areas of density across the dataset (e.g. DBSCAN⁵) or relationships to statistical distributions⁶.

In many cases the number of groups that should be present in a dataset is not known in advance and deciding the correct number of clusters to use is a significant challenge. For some algorithms, such as k -means clustering, the number of clusters must be explicitly provided. Other methods have parameters that, directly or

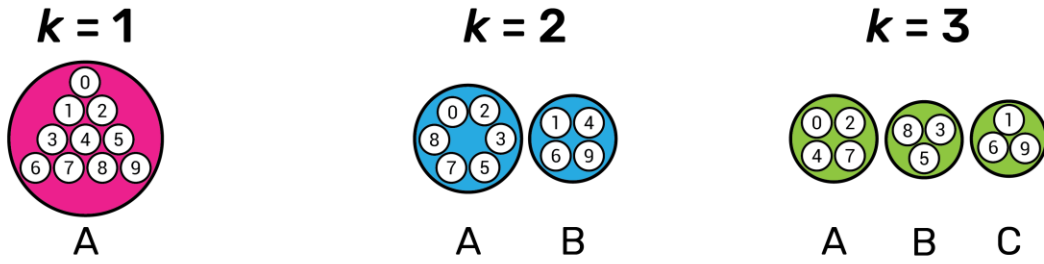
1 indirectly, control the clustering resolution and therefore the number of clusters
2 produced. While there are methods and statistics (such as the elbow method⁷ or
3 silhouette plots⁸) designed to help analysts decide which clustering resolution to use,
4 they typically produce a single score which only considers a single set of samples or
5 clusters at a time.
6
7

8
9 An alternative approach would be to consider clusterings at multiple resolutions and
10 examine how samples change groupings as the number of clusters increases. This is
11 the approach taken by the clustering tree visualisation we present here: (i) a dataset
12 is clustered at multiple resolutions producing sets of cluster nodes, (ii) the overlap
13 between clusters at adjacent resolutions is used to build edges, (iii) the resulting
14 graph is presented as a tree. This tree can be used to examine how clusters are related
15 to each other, which clusters are distinct and which are unstable. In the following
16 sections we describe how we construct such a tree and present examples of trees built
17 from a classical clustering dataset and a complex single-cell RNA-sequencing
18 (scRNA-seq) dataset. The figures shown here can be produced in R using our publicly
19 available clustree package.
20
21
22
23
24
25
26
27
28

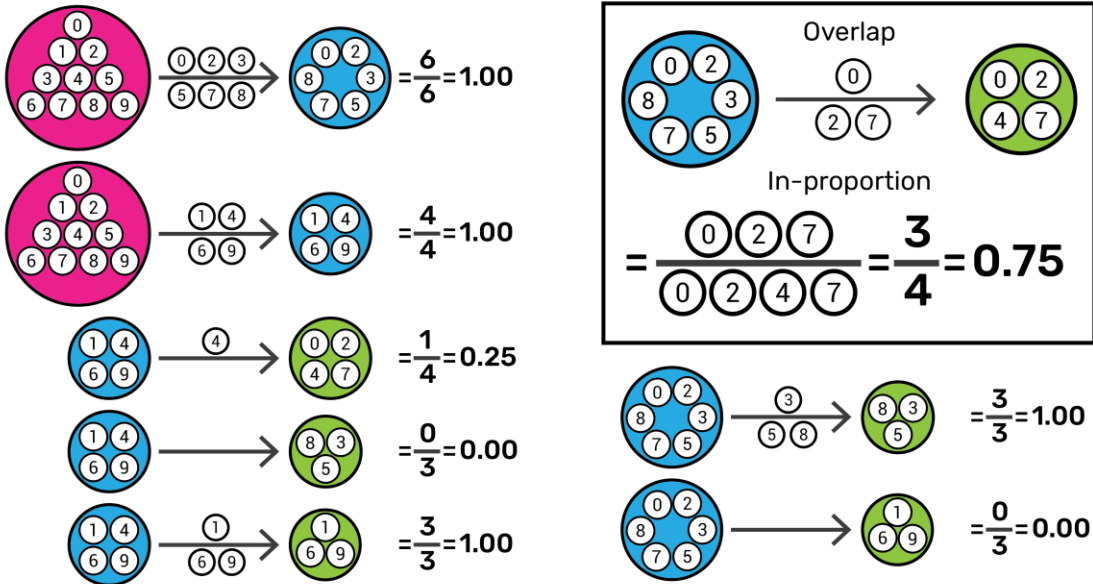
29 Building a clustering tree

30
31 To build a clustering tree, we start with a set of clusterings allocating samples to
32 groups at several different resolutions. These could be produced using any hard-
33 clustering algorithm that allows control of the number of clusters in some way. For
34 example, this could be a set of samples clustered using k -means with $k = 1, 2, 3$ as
35 shown in Figure 1. We sort these clusterings so that they are ordered by increasing
36 resolution (k), then consider pairs of adjacent clusterings. Each cluster $c_{k,i}$ (where $i =$
37 $1, \dots, n$ and n is the number of clusters at resolution k) is compared with each cluster
38 $c_{k+1,j}$ (where $j = 1, \dots, m$ and m is the number of clusters at resolution $k + 1$). The
39 overlap between the two clusters is computed as the number of samples that are
40 assigned to both $c_{k,i}$ and $c_{k+1,j}$. We next build a graph where each node is a cluster
41 and each edge is an overlap between two clusters.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1. Cluster at multiple resolutions



2. Find overlaps and calculate in-proportion



3. Filter edges and visualise tree

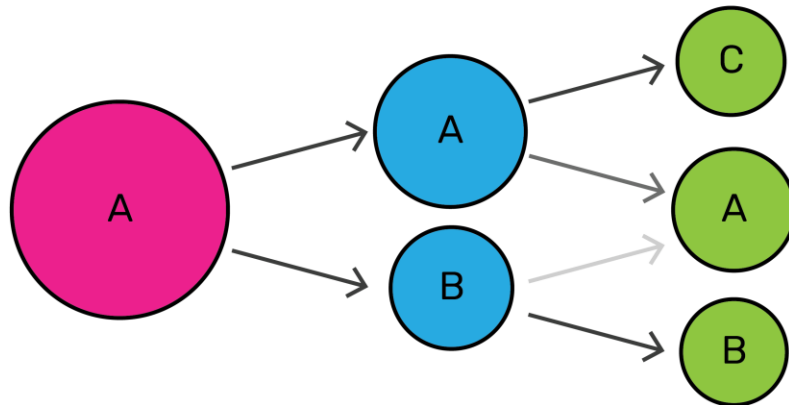


Figure 1 Illustration of the steps required to build a clustering tree. First a dataset must be clustered at different resolutions. The overlap in samples between clusters at adjacent resolutions is computed and used to calculate the in-proportion for each edge. Finally the edges are filtered and the graph visualised as a tree.

Many of the edges will be empty, for example in Figure 1 no samples in Cluster A at $k = 2$ end up in Cluster B at $k = 3$. In some datasets there may also be edges that contain few samples. These edges are not informative and result in a cluttered tree. An obvious solution for removing uninformative, low-count edges is to filter them

1 using a threshold on the number of samples they represent. However, in this case the
2 count of samples is not the correct statistic to use because it favours edges at lower
3 resolutions and those connecting larger clusters. Instead we define the in-proportion
4 metric as the ratio between the number of samples on the edge and the number of
5 samples in the cluster it goes towards. This metric shows the importance of the edge
6 to the higher resolution cluster independently of the cluster size. We apply a
7 threshold to the in-proportion in order to remove less informative edges.
8
9

10
11
12 The final graph is visualised using a tree layout. This places the cluster nodes in a
13 series of layers where each layer is a different clustering resolution and edges show
14 the transition of samples through those resolutions. Edges are coloured according to
15 the number of samples they represent and the in-proportion metric is used to control
16 the edge transparency, highlighting more important edges. By default, the size of
17 nodes is adjusted according to the number of samples in the cluster and their colour
18 indicates the resolution. The clustree package also includes options for controlling
19 the aesthetics of nodes based on the attributes of samples in the clusters they
20 represent.
21
22
23
24
25
26
27
28

29 **A simple example**

30
31 To further illustrate how a clustering tree is built, we will work through an example
32 using the classical iris dataset⁹. This dataset contains measurements of the sepal
33 length, sepal width, petal length and petal width from 150 iris flowers, 50 from each
34 of three species: *Iris setosa*, *Iris versicolor* and *Iris virginica*. The iris dataset is
35 commonly used as example for both clustering and classification problems with the
36 *Iris setosa* samples being significantly different to, and linearly separable from, the
37 other samples. We have clustered this dataset using k -means clustering with $k =$
38 $1, \dots, 5$ and produced the clustering tree shown in Figure 2A.
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

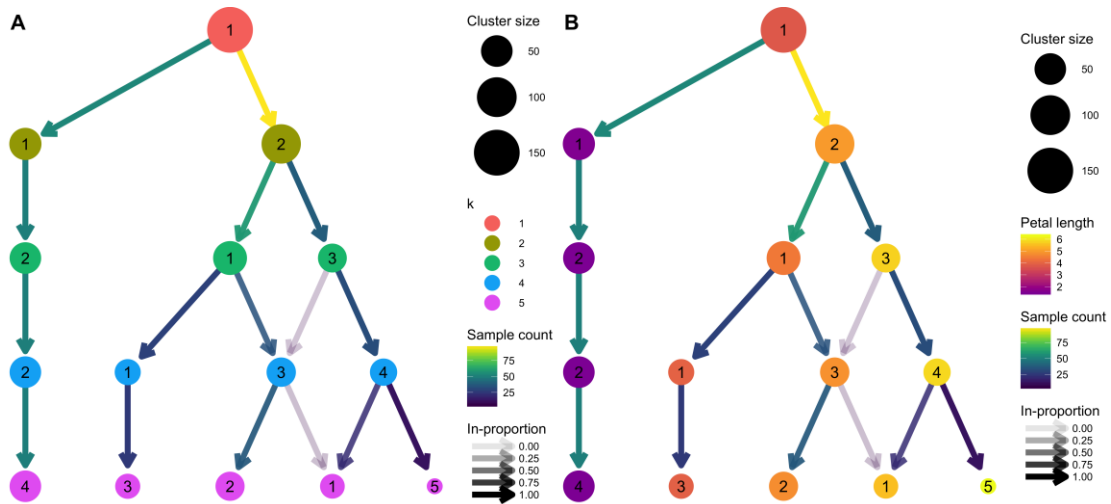


Figure 2 Clustering trees based on k -means clustering of the iris dataset. In A, nodes are coloured according to the value of k and sized according to the number of samples they represent. Edges are coloured according to the number of samples (from blue representing few to yellow representing many) and the transparency adjusted according to the in-proportion, with stronger lines showing edges that are more important to the higher resolution cluster. Cluster labels are randomly assigned by the k -means algorithm. B shows the same tree with the node colouring changed to show the mean petal length of the samples in each cluster.

We see that there is one branch of the tree that is clearly distinct (presumably representing *Iris setosa*), remaining unchanged regardless of the number of clusters. On the other side we see the cluster at $k = 2$ cleanly split into two clusters (presumably *Iris versicolor* and *Iris virginica*) at $k = 3$ but as we move to $k = 4$ and $k = 5$ we see clusters being formed from multiple branches with more low proportion edges. This kind of pattern indicates that the data has become over-clustered and we have begun to introduce artificial groupings. In this case we know that $k = 3$ is the correct choice but this is also the value that is suggested by this tree.

We can check our assumption that the distinct branch represents the *Iris setosa* samples and the other two clusters at $k = 3$ are *Iris versicolor* and *Iris virginica* by overlaying some known information about the samples. In Figure 2B we have coloured the nodes by the mean petal length of the samples they contain. We can now see that clusters in the distinct branch have the shortest petals, with Cluster 1 at $k = 3$ having an intermediate length and Cluster 3 the longest petals. This feature is known to separate the samples into the expected species with *Iris setosa* having the shortest petals on average, *Iris versicolor* an intermediate length and *Iris virginica* the longest.

Although this is a very simple example it still highlights some of the benefits of viewing a clustering tree. We get some indication of the correct clustering resolution

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

by examining the edges and we can overlay known information to assess the quality of the clustering. For example, if we observed that all clusters had the same mean petal length it would suggest that the clustering has not been successful as we know this is an important feature that separates the species. We could potentially learn more by looking at which samples follow low proportion edges or overlaying a series of features to try and understand what causes particular clusters to split.

Clustering trees for single-cell RNA-seq data

One field that has begun to make heavy use of clustering techniques is the analysis of single-cell RNA-sequencing (scRNA-seq) data. Single-cell RNA-sequencing is a recently developed technology that can measure how genes are expressed in thousands to millions of individual cells¹¹. This technology has been rapidly adopted in fields like developmental biology and immunology where it is valuable to have information from single cells rather than measurements that are averaged across the many different cells in a sample using older RNA sequencing technologies. One of the key uses for scRNA-seq is to discover and interrogate the different cell types present in a sample of a complex tissue. In this situation, clustering is typically used to group similar cells based on their gene expression profiles. Differences in gene expression between groups can then be used to infer the identity or function of those cells¹². The number of cell types in an scRNA-seq dataset can vary depending on factors such as the tissue being studied, its developmental or environmental state and the number of cells captured. Often the number of cells types is not known before the data is generated and some samples can contain dozens of clusters. Therefore, deciding which clustering resolution to use is an important consideration in this application.

As an example of how clustering trees can be used in the scRNA-seq context we consider a commonly used Peripheral Blood Mononuclear Cell (PBMC) dataset. This dataset was originally produced by 10x Genomics and contains 2700 peripheral blood mononuclear cells, representing a range of well-studied immune cell types¹³. We have analysed this dataset using the Seurat package¹⁴, a commonly used toolkit for scRNA-seq analysis, following the instructions in their tutorial with the exception of varying the clustering resolution parameter from zero to five (see methods). Seurat uses a graph-based clustering algorithm and the resolution parameter controls the partitioning of this graph, with higher values resulting in more clusters. The clustering trees produced from this analysis are shown in Figure 3.

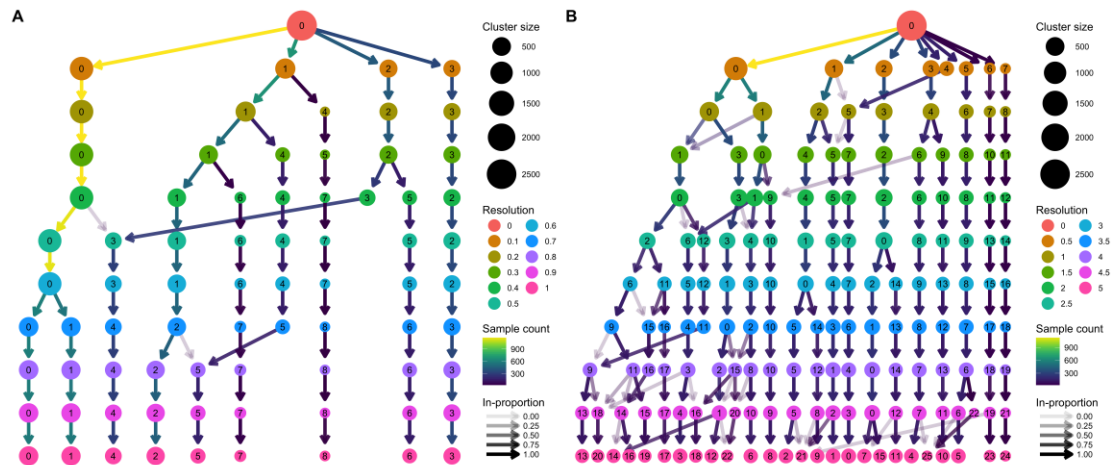


Figure 3 Two clustering trees of a dataset of 2700 Peripheral Blood Mononuclear Cells (PBMCs). A) results from clustering using Seurat with resolution parameters from zero to one. At a resolution of 0.1 we see the formation of four main branches, one of which continues to split up to a resolution of 0.5, after which there are only minor changes. B) resolutions from zero to five. At the highest resolutions we begin to see many low in-proportion edges indicating cluster instability. Seurat labels clusters according to their size with Cluster 0 being the largest.

The clustering tree covering resolutions zero to one in steps of 0.1 (Figure 3A) shows that four main branches form at a resolution of just 0.1. One of these branches, starting with Cluster 3 at resolution 0.1, remains unchanged while the branch starting with Cluster 2 splits only once at a resolution of 0.4. Most of the branching occurs in the branch starting with Cluster 1 which consistently has sub-branches split off to form new clusters as the resolution increases. There are two regions of stability in this tree; at resolution 0.5-0.6 and resolution 0.7-1.0 where the branch starting at Cluster 0 splits in two.

Figure 3B shows a clustering tree with a greater range of resolutions, from zero to five in steps of 0.5. By looking across this range we can see what happens when the algorithm is forced to produce more clusters than are likely to be truly present in this dataset. As over-clustering occurs we begin to see more low in-proportion edges and new clusters forming from multiple parent clusters. This suggests that those areas of the tree are unstable and that the new clusters being formed are unlikely to represent true groups in the dataset.

Known marker genes are commonly used to identify the cell types that specific clusters correspond to. Overlaying gene expression information onto a clustering tree provides an alternative view that can help to indicate when clusters containing pure cell populations are formed. Figure 4 shows the PBMC clustering tree in Figure 3A overlaid with the expression of some known marker genes.

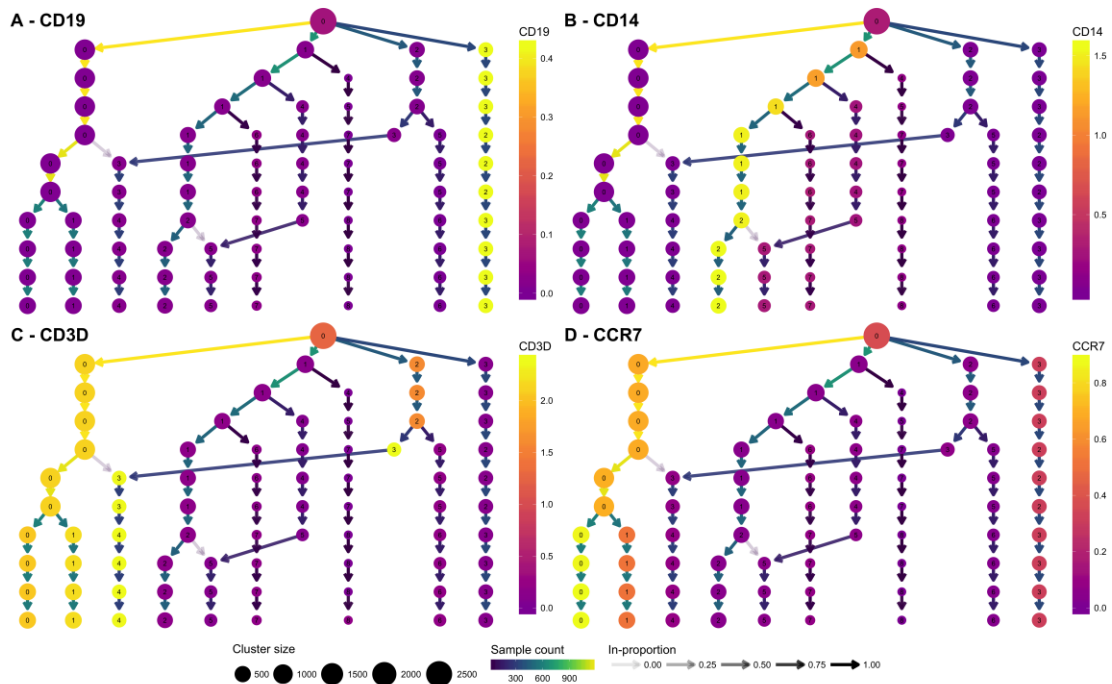


Figure 4 Clustering trees of the PBMC dataset coloured according to the expression of known markers. The node colours indicate the average of the \log_2 gene counts of samples in each cluster. CD19 (A) identifies B cells, CD14 (B) shows a population of monocytes, CD3D (C) is a marker of T cells and CCR7 (D) shows the split between memory and naive CD4 T cells.

By adding this extra information, we can quickly identify some of the cell types. CD19 (Figure 4A) is a marker of B cells and is clearly expressed in the most distinct branch of the tree. CD14 (Figure 4B) is a marker of a type of monocyte, which becomes more expressed as we follow one of the central branches, allowing us to see which resolution identifies a pure population of these cells. CD3D (Figure 4C) is a general marker of T cells and is expressed in two separate branches, one which splits into low and high expression of CCR7 (Figure 4D), separating memory and naive CD4 T cells. By adding expression of known genes to a clustering tree, we can see if more populations can be identified as the clustering resolution is increased and if clusters are consistent with known biology. For most of the Seurat tutorial a resolution of 0.6 is used, but the authors note that by moving to a resolution of 0.8, a split can be achieved between memory and naive CD4 T cells. This is a split that could be anticipated by looking at the clustering tree.

Discussion and conclusion

Clustering similar samples into groups is a useful technique in many fields, but often analysts are faced with the tricky problem of deciding which clustering resolution to use. Traditional approaches to this problem typically consider a single cluster or

1 sample at a time and may rely on prior knowledge of sample labels. Here we present
2 clustering trees, an alternative visualisation that shows the relationships between
3 clusterings at multiple resolutions.
4

5 Clustering trees display how clusters are divided as resolution increases, which
6 clusters are clearly separate and distinct, which are related to each other and how
7 samples change groups as more clusters are produced. Although clustering trees can
8 appear similar to the trees produced from hierarchical clustering there are several
9 important differences. Hierarchical clustering considers the relationships between
10 individual samples and doesn't provide an obvious way to form groups. In contrast,
11 clustering trees are independent of any particular clustering method and show the
12 relationships between distinct groups of samples, any of which could be used for
13 further analysis.
14
15
16
17
18
19
20

21 To illustrate the uses of clustering trees we presented two examples, one using the
22 classical iris dataset and a second based on a complex scRNA-seq dataset. Both
23 examples demonstrate how a clustering tree can suggest the correct resolution to use
24 and how overlaying extra information can help to validate those clusters. This is of
25 particular use to scRNA-seq analysis as these datasets are often large, noisy and
26 contain an unknown number of cell types.
27
28
29
30

31 Even when the number of clusters to choose is not a problem, clustering trees can be
32 a valuable tool. They provide a compact, information dense, visualisation that can
33 display summarised information across a range of clusters. By modifying the
34 appearance of cluster nodes based on attributes of the samples they represent,
35 clusterings can be evaluated and identities of clusters established. Clustering trees
36 potentially have applications in many fields and in the future could be adapted to be
37 more flexible, such as by accommodating fuzzy clusterings.
38
39
40
41
42
43
44
45

46 **Methods**

47 The clustree software package is built for the R statistical programming language. It
48 relies on the ggraph package (<https://github.com/thomasp85/ggraph>), which is itself
49 built on the ggplot2¹⁵ and tidygraph packages
50 (<https://github.com/thomasp85/tidygraph>). Clustering trees are displayed using the
51 Reingold-Tilford tree layout¹⁶ or the Sugiyama layout¹⁷, both available as part of the
52 igraph package¹⁸.
53
54
55
56
57
58
59
60
61
62
63
64
65

1 The iris dataset is available as part of R. We clustered this dataset using the “kmeans”
2 function in the stats package with values of k from one to five. Each value of k was
3 clustered with a maximum of 100 iterations and with 10 random starting positions.
4 The clustered iris dataset is available as part of the clustree package.
5
6

7 The PBMC dataset was downloaded from the Seurat tutorial page
8 (http://satijalab.org/seurat/pbmc3k_tutorial.html) and this tutorial was followed for
9 most of the analysis. Briefly cells were filtered based on the number of genes they
10 express and the percentage of counts assigned to mitochondrial genes. The data was
11 then log-normalised and 1838 variable genes identified. Potential confounding
12 variables (number of unique molecular identifiers and percentage mitochondrial
13 expression) were regressed from the dataset before performing principal component
14 analysis on the identified variable genes. The first 10 principal components were then
15 used to build a graph which was partitioned into clusters using Louvain modularity
16 optimisation¹⁹ with resolution parameters in the range zero to five, in steps of 0.1
17 between zero and one and then in steps of 0.5.
18
19
20
21
22
23
24
25
26

27 **Declarations**

28 **Ethics**

29 Not applicable.
30
31
32
33
34

35 **Availability of data and materials**

36 The clustree package is available from GitHub at <https://github.com/lazappi/clustree>
37 and the code and datasets used for the analysis in this paper are available from
38 <https://github.com/Oshlack/clustree-paper>. The clustered iris dataset is included as
39 part of clustree and the PBMC dataset can be downloaded from the Seurat tutorial
40 page (http://satijalab.org/seurat/pbmc3k_tutorial.html) or the paper GitHub
41 repository.
42
43
44
45
46
47

48 **Competing interests**

49 The authors declare no competing interests.
50
51
52

53 **Funding**

54 Luke Zappia is supported by an Australian Government Research Training Program
55 (RTP) Scholarship. Alicia Oshlack is supported through a National Health and
56 Medical Research Council Career Development Fellowship APP1126157. MCRI is
57
58
59
60
61
62
63
64
65

supported by the Victorian Government's Operational Infrastructure Support Program.

Acknowledgements

Thank you to Marek Cmero for providing comments on a draft of the manuscript.

References

1. Forgy, W. E. Cluster analysis of multivariate data : Efficiency versus interpretability of classifications. *Biometrics* **21**, 768–769 (1965).
2. Macqueen, J. Some methods for classification and analysis of multivariate observations. in *In 5th berkeley symposium on mathematical statistics and probability* (1967).
3. Lloyd, S. Least squares quantization in PCM. *IEEE Trans. Inf. Theory* **28**, 129–137 (1982).
4. Kaufman, L. & Rousseeuw, P. J. Partitioning around medoids (program PAM). in *Finding groups in data* 68–125 (John Wiley & Sons, Inc., 1990).
5. Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. A density-based algorithm for discovering clusters a density-based algorithm for discovering clusters in large spatial databases with noise. in *Proceedings of the second international conference on knowledge discovery and data mining* 226–231 (AAAI Press, 1996).
6. Fraley, C. & Raftery, A. E. Model-Based clustering, discriminant analysis, and density estimation. *J. Am. Stat. Assoc.* **97**, 611–631 (2002).
7. Thorndike, R. L. Who belongs in the family? *Psychometrika* **18**, 267–276 (1953).
8. Rousseeuw, P. J. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.* **20**, 53–65 (1987).
9. Anderson, E. The Irises of the Gaspé Peninsula. *Bulletin of the American Iris Society* **59**, 2–5 (1935).
10. Fisher, R. A. The use of multiple measurements in taxonomic problems. *Ann. Eugen.* **7**, 179–188 (1936).
11. Tang, F. *et al.* mRNA-Seq whole-transcriptome analysis of a single cell. *Nat. Methods* **6**, 377–382 (2009).
12. Stegle, O., Teichmann, S. A. & Marioni, J. C. Computational and analytical challenges in single-cell transcriptomics. *Nat. Rev. Genet.* **16**, 133–145 (2015).
13. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* **8**, 14049 (2017).
14. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33**, 495–502 (2015).
15. Wickham, H. *Ggplot2: Elegant graphics for data analysis*. (Springer New York, 2010).

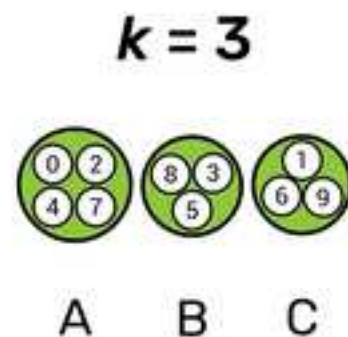
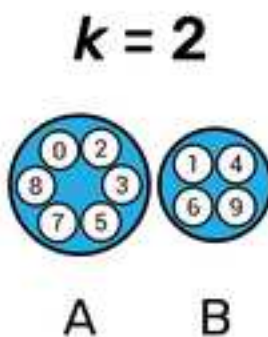
16. Reingold, E. M. & Tilford, J. S. Tidier drawings of trees. *IEEE Trans. Software Eng.* **SE-7**, 223–228 (1981).

17. Sugiyama, K., Tagawa, S. & Toda, M. Methods for visual understanding of hierarchical system structures. *IEEE Trans. Syst. Man Cybern.* **11**, 109–125 (1981).

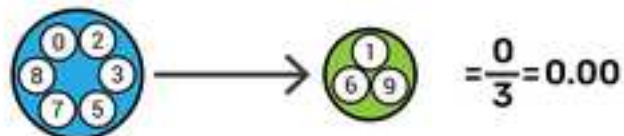
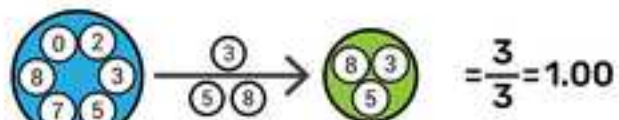
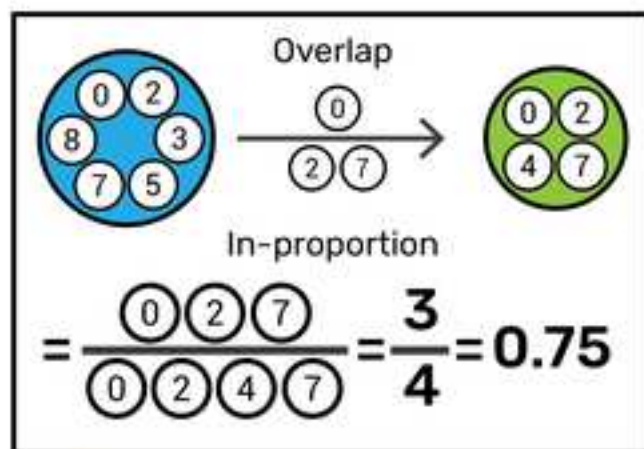
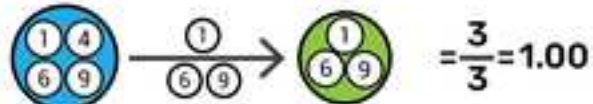
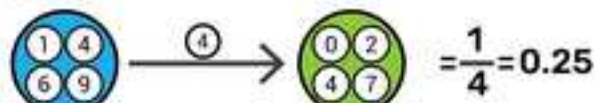
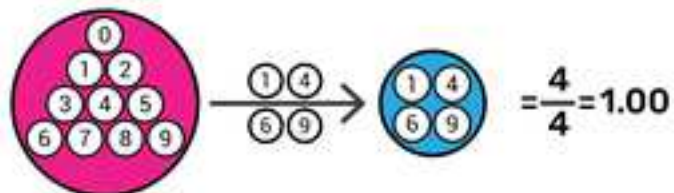
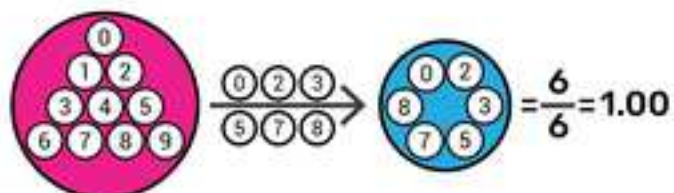
18. Csardi, G. & Nepusz, T. The igraph software package for complex network research. *InterJournal, Complex Systems* **1695**, 1–9 (2006).

19. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat. Mech.* **2008**, P10008 (2008).

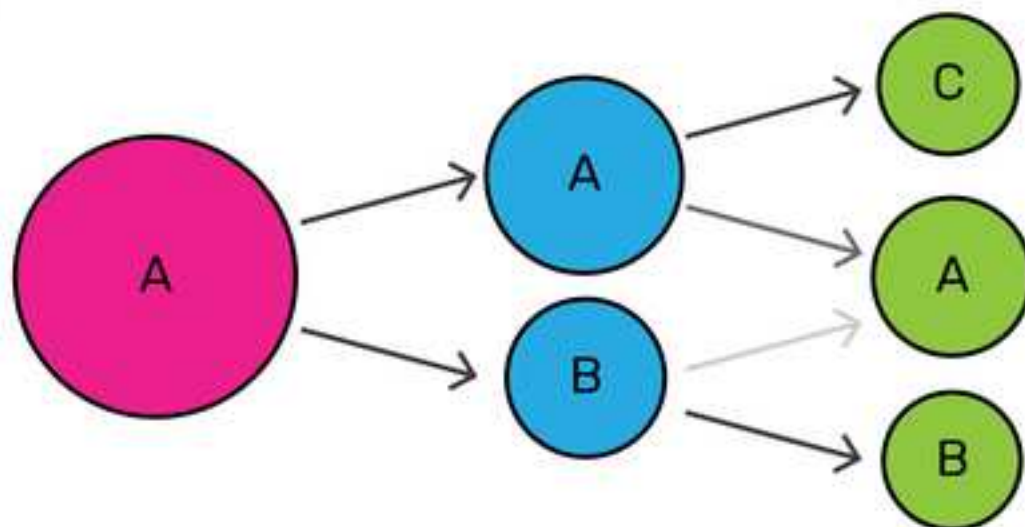
1. Cluster at multiple resolutions



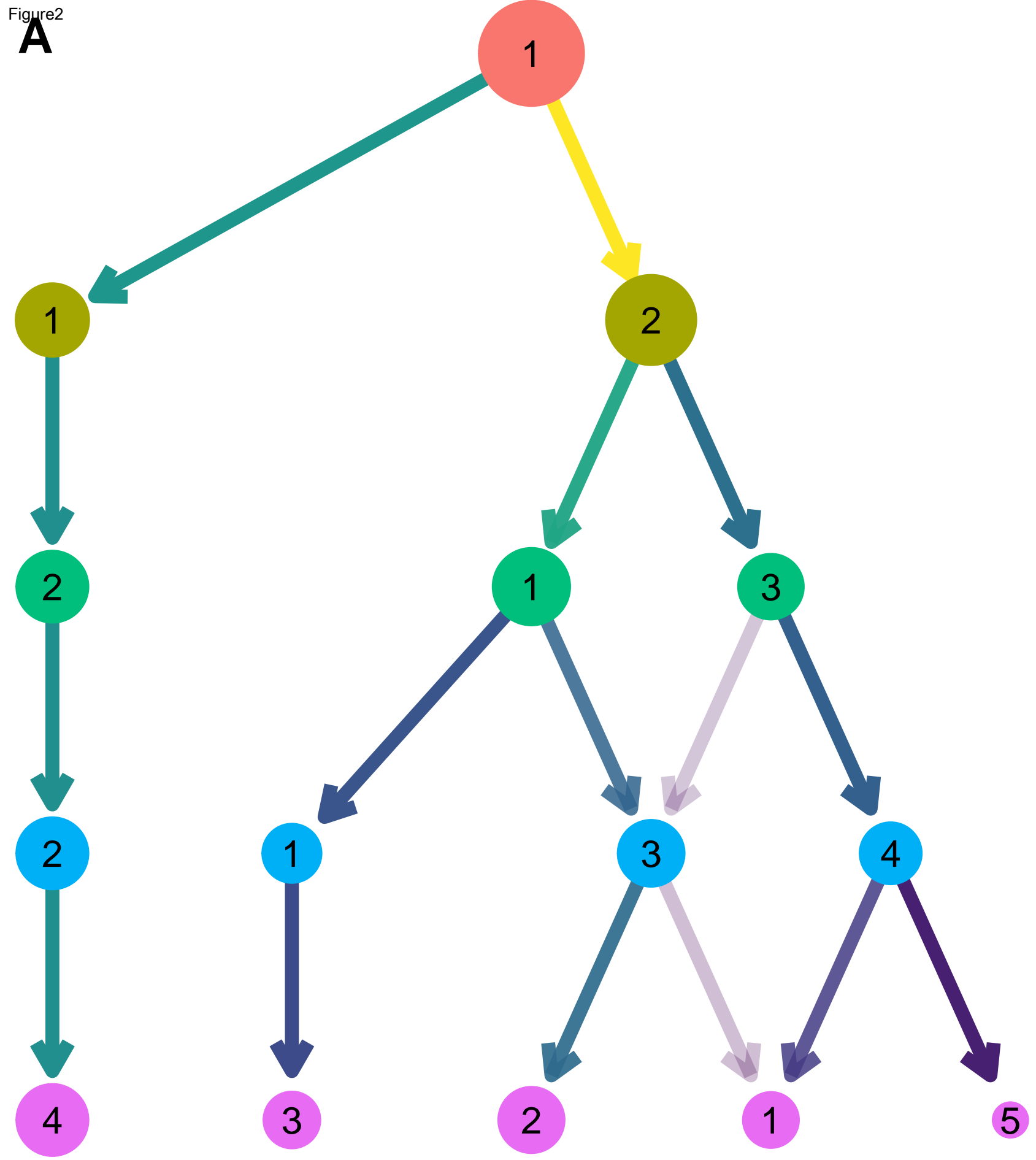
2. Find overlaps and calculate in-proportion



3. Filter edges and visualise tree



A



B

