


Dimensionality reduction by UMAP to visualize physical and genetic interactions

Michael W. Dorrity ^{1,3}, Lauren M. Saunders^{1,3}, Christine Queitsch¹, Stanley Fields ^{1,2✉} & Cole Trapnell ^{1✉}

Dimensionality reduction is often used to visualize complex expression profiling data. Here, we use the Uniform Manifold Approximation and Projection (UMAP) method on published transcript profiles of 1484 single gene deletions of *Saccharomyces cerevisiae*. Proximity in low-dimensional UMAP space identifies groups of genes that correspond to protein complexes and pathways, and finds novel protein interactions, even within well-characterized complexes. This approach is more sensitive than previous methods and should be broadly useful as additional transcriptome datasets become available for other organisms.

¹Department of Genome Sciences, University of Washington, Seattle, WA 98195, USA. ²Department of Medicine, University of Washington, Seattle, WA 98195, USA. ³These authors contributed equally: Michael W. Dorrity, Lauren M. Saunders. ✉email: fields@uw.edu; coletrap@uw.edu

A central goal of biological studies is the identification and characterization of proteins that act in a common cellular pathway. Efforts toward this goal have been greatly aided by large-scale perturbation analyses coupled with whole-transcriptome profiling, in which each gene's transcriptional response to a perturbation is measured. If a sufficient database of expression profiles exists, then a pathway affected by an uncharacterized perturbation such as a gene mutation, drug treatment or growth condition—can be described by matching the resultant profile to a known profile¹. For the yeast *Saccharomyces cerevisiae*, the expression profiles of a large number of individual yeast deletion mutants have been established and used to infer protein complexes and networks^{2–4}. Maximizing the utility of expression profiling approaches for inference of physical and genetic interactions requires ever larger such datasets. However, standard techniques, such as pairwise correlation, do not fully capture the variation available to link gene function as more dimensions are added from larger scale experiments. Therefore, techniques that reduce dimensionality of the data while maintaining relationships between genes are imperative for the inference of physical and genetic interactions in very large gene expression datasets.

Dimensionality reduction methods capture variability in a limited number of random variables to facilitate 2- or 3D-visualization of datasets with tens to thousands of dimensions. This approach is recognizable in the commonly used method of principal component analysis (PCA), which uses linear combinations of variables to generate orthogonal axes that efficiently capture the variation present in the data with fewer variables. Another approach, t-Distributed Stochastic Neighbor Embedding (t-SNE), carries out dimensionality reduction by analyzing similarity of points using a Gaussian distance in high-dimensional space and projecting these data into a low-dimensional space⁵. A more recent method, Uniform Manifold Approximation and Projection (UMAP), estimates a topology of the high-dimensional data and uses this information to construct a low-dimensional representation that preserves relationships present in the data⁶. UMAP has been particularly useful to precisely define cell types in mixed populations based on data from single-cell RNA-seq experiments^{7–13}; it also performs well on other gold-standard datasets^{6,14}. Because UMAP is better able to preserve elements of the data structure from high-dimensional space than similar outputs from t-SNE, it captures local relationships within groups of transcriptomes in addition to global relationships between distinct groups¹⁴. This feature is especially useful in the inference of gene relationships, which can be due to physical interaction, overlapping gene function, or coordinated contributions to a larger cellular process. Here, we show that the use of dimensionality reduction by UMAP on bulk expression profiling data of 1484 single-gene mutants of *S. cerevisiae* links gene function in clusters at increasingly finer scales, corresponding to broad cellular activities, pathways, protein complexes and individual protein-protein interactions.

Results

UMAP groups deletion mutants with shared protein function.

We assigned groups, or clusters, to deletion mutants with similar transcriptional responses using the Louvain community detection algorithm in low-dimensional UMAP space⁹. While many single-cell transcriptomic studies use expression values from genes with the highest dispersion across individual cells, we took advantage of the completeness of bulk microarray data generated by Kemmeren et al.³ and used expression values for all 6170 genes measured in each of the 1484 single-gene deletion strains to make a UMAP projection for subsequent clustering. This approach resolved 50 main clusters, with the number of deletion backgrounds assigned to each cluster ranging from 4 to 298 (median of 11). Clusters with >25 strains were subsequently sub-clustered using similar

parameters to define groups. The final dataset contains 171 clusters with a median of 8 strains per cluster.

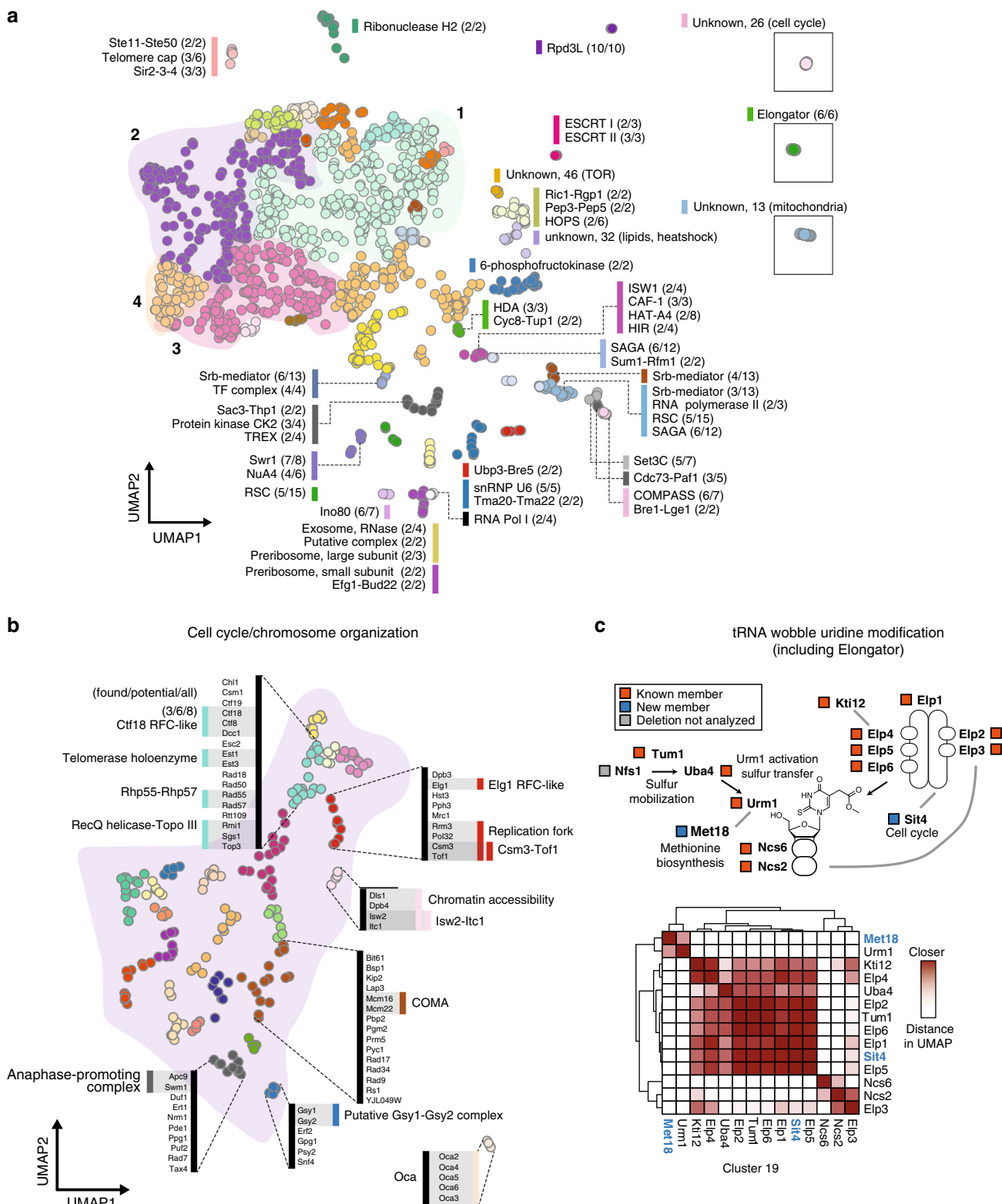
A total of 194 characterized yeast complexes have at least two of their corresponding genes in the dataset of single deletions. For 40% of these complexes (78/194), we could assign two or more genes to the same cluster (examples of complexes in the initial set of 50 clusters in Fig. 1a, additional complexes were separated in the sub-clustered set (Fig. 1b)). For example, the sub-clustering of the original cluster 2, which is characterized by cell cycle and chromosome organization genes, resulted in more distinctly separating the Isw2-Itc1 chromatin remodeling complex, the Csm3-Tof1 S-phase checkpoint complex and the Oca S-phase histone activation complex (Fig. 1b). Within this sub-clustered set, multiple complexes could be found among genes within a single cluster, suggesting that these complexes may cooperatively contribute to chromosome cohesion and recombination (Fig. 1b).

In some cases, members of individual complexes were assigned to separate clusters, suggesting sub-functionalization of components. For example, the 13-member mediator complex was found in three clusters (numbers 16, 34, and 41) containing 3, 6, and 4 members of mediator, respectively (Fig. 1a). Cluster 16 also contains members of SAGA and SWI/SNF complexes, and loss of mediator subunits in this cluster alters the transcription of amino acid metabolism genes and glucose transmembrane transporters (Supplementary Data 1); cluster 34 contains galactose-responsive subunits of mediator; and cluster 41 contains transcriptional initiation-related mediator subunits. Here, UMAP preserves global relationships between clusters in addition to resolving proximal cluster members. For example, most chromatin remodeling complexes grouped in UMAP space, despite being present in separate clusters and containing unique local topologies (Fig. 1a).

UMAP clustering identified the components of the pathway for tRNA wobble uridine modification (Fig. 1c), which requires the *URM1* pathway for 2-thiolation and the Elongator complex for side chain formation at U₃₄ of tRNA¹⁵. The clustering revealed two additional members that are likely to link metabolism and cell cycle to this process. One of these, Met18, has a human ortholog (MMS19) that functions in maturation of Fe-S cluster-containing proteins; the conserved yeast and human Elongator component Elp3 is one of these Fe-S proteins¹⁶. The other new member, the PP2A phosphatase Sit4, is implicated in dephosphorylation of Elongator; its absence leads to tRNA modification defects¹⁷.

Comparison of UMAP distance to other protein interaction metrics.

To assess whether distance in UMAP space captured known interactions as well as pairwise correlation, we used three gold-standard interaction datasets: (1) protein interactions determined by co-immunoprecipitation followed by mass spectrometry²; (2) gene interactions from stringDB¹⁸, which are derived from a probabilistic metric based on multiple evidence channels including yeast two-hybrid, pathway annotations, and co-expression; and (3) interactions from CellMAP¹⁹, which are derived from an experimental screen for synthetic genetic interactions. The UMAP distance metric captured protein complexes more sensitively and with more precision than previous pairwise correlation-based metrics (AUC pairwise correlation = 0.73, AUC UMAP = 0.84, Fig. 2a). UMAP distance also captured known interacting pairs better than distance in high-dimensional space (AUC = 0.56) and distance in PCA space (AUC = 0.70), suggesting that the UMAP dimensionality reduction itself adds value in the identification of interactions (Fig. 2a, Supplementary Fig. 1a). Across each gold-standard interaction dataset, UMAP distance performed better than several other standard approaches for analyzing the transcriptome data, including PCA, random orthonormal projections²⁰, and tSNE (Supplementary Fig. 1a, b).



Performing clustering in UMAP space ought to produce clusters containing more true interactions than distance in other spaces. To test whether similar results were obtained without UMAP dimensionality reduction, we clustered the data in PCA space. Clustering in PCA space identified 8/50 clusters with perfect overlap to UMAP clusters, and 34/50 that overlap by at least 50% (Supplementary Fig. 1c).

To compare pairwise correlation with the UMAP approach, we calculated for each known interacting pair (1) the Pearson correlation of their deletion transcriptomes; and (2) the distance of those two genes in the UMAP space generated by using all deletion transcriptomes. Among these interacting pairs, UMAP distance and pairwise correlation are negatively correlated (Fig. 2b). However, the increased sensitivity of UMAP distance

Fig. 1 UMAP clusters single-gene deletion transcriptomes according to shared function. **a** UMAP coordinates of 1484 single-gene deletion strains clustered by similarity in transcriptional effects. The initial 50 individual clusters are each shown in a different color. Strains that comprise protein complexes are indicated alongside a bar colored according to cluster identity. Each complex is represented as a fraction: the number of complex members found in the cluster over the number of complex members in the set of 1484 mutants. Clusters with coordinates far from the main group are shown in boxes. Clusters without a known complex are marked as “unknown,” along with an arbitrary cluster number; these clusters are annotated with a broad GO term enriched in that cluster. **b** Cluster 2 shows more distinct groupings when re-clustered separately. Annotations as in **a**. Cluster 2 as a whole was enriched for cell cycle and chromosome organization, with individual clusters corresponding to parts of this process. **c** The tRNA wobble uridine pathway, captured entirely within the cluster containing the Elongator complex (boxed green cluster in **a**). Complex members within this cluster are annotated with orange boxes, while new members are annotated in blue. One pathway member, *Nfs1*, was not present in the single-gene deletion dataset. The heatmap represents fine-scale distances between each pair of points within the cluster. Darker shades of red indicate points nearer in UMAP space; hierarchical clustering was applied on this distance metric to group proteins within this pathway. Heterodimeric interactions, such as *Ncs6-Ncs2* (bottom-right corner of heatmap), are nearer to each other than other members of the pathway. Novel members of this pathway (blue text) are grouped with other members based on their similarity of UMAP distance, and these new interactions are indicated with gray lines in the pathway diagram.

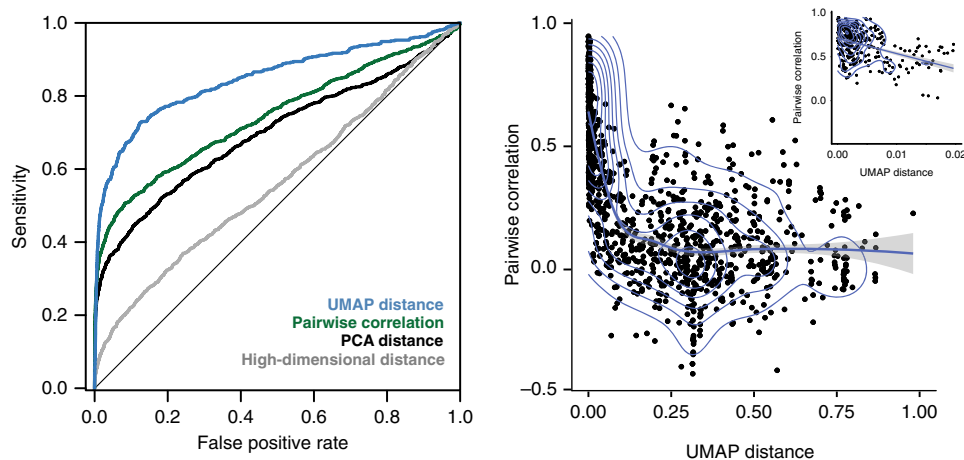


Fig. 2 UMAP distance identifies protein-protein interactions more effectively than previous methods. **a** A receiver-operator curve showing the ability of UMAP distance to capture known protein-protein interactions (sensitivity) as a function of its false positive detection. UMAP distance (blue) performs better than pairwise correlation (green), PCA distance (dark gray), and high-dimensional distance (light gray) in identifying interactions. **b** For each protein-protein interaction, the distance between points in UMAP space was plotted against the pairwise correlation of that pair of transcriptomes. The density of points is indicated with blue lines. Inset in the upper right shows a zoomed-in portion of the x-axis; points with UMAP distance in this range are highly enriched for true interactions that are not captured by pairwise correlation.

to detect known interactions suggests that the discrepancies between UMAP distance and pairwise correlation might represent interactions that were previously overlooked. Based on a UMAP distance cutoff corresponding to a 5% FDR of known complex members (Inset, Fig. 2b), we were able to identify 176 putative interactions that would not have been confidently called by previous approaches using pairwise correlations ($PCC < 0.5$); these interactions contain 86 unique genes, of which 77 show co-IP or yeast two-hybrid evidence for membership among 31 protein complexes, while the remaining 9 genes had no such evidence.

Since proximity in UMAP space tends to capture known interactions and shared function, distance in UMAP space could serve as a useful tool to investigate evolutionary questions about gene divergence. We calculated UMAP distance between 151 paralogous gene pairs in yeast and used this distance to characterize the functional divergence between each pair (Supplementary Fig. 2a). Proximity of paralog pairs in UMAP space did not correspond to previous estimations of paralog divergence (Supplementary Fig. 2b, c) based on synthetic genetic interaction ($R = 0.018$) or Gene Ontology relationships ($R = 0.035$)¹⁹. When paralogs show a negative genetic interaction—that is, deletion of both genes leads to lower fitness than expected—it is assumed that the two genes retain redundant functions. However, 11 paralog pairs whose negative genetic interactions suggested redundant function showed distinct downstream effects on gene expression when each gene was

deleted (Supplementary Fig. 2b, d); these genes may have distinct effects on fitness in different environments²¹. In these cases, a gene may retain the capacity to complement the essential function of its paralogous partner, while diverging sufficiently in function as revealed by the UMAP-based transcriptome analysis.

Despite successful clustering of many protein complexes and pathways of yeast, the UMAP approach nevertheless identified several clusters that did not obviously correspond to a complex or pathway. We used GO enrichment of differentially expressed genes in these clusters to interrogate their function: cluster 26 showed enriched terms for cell cycle, non-membrane-bound organelles, and prions; cluster 13 showed enrichment for mitochondrial function; cluster 46 showed enrichment for TOR signaling and aerobic respiration; cluster 32 showed enrichment for protein folding; and cluster 11 showed enrichment for heme binding. Differential expression analysis produced significant gene sets for all main and sub-clusters (Supplementary Data 1).

Discussion

Because of its greater sensitivity than other approaches, as well as its ability to capture both local and global relationships, UMAP-based association of gene function adds value in the identification of protein complexes, pathways, and novel interactions in transcriptomic datasets. However, the utility of this method is

dependent on the availability of high-quality profiling data from large-scale environmental or genetic perturbation experiments. As more datasets of this type become available, we expect that this approach, or similar dimensionality reduction techniques, will become increasingly useful in mapping protein complexes and pathways both within and across other species. The recent appearance of single-cell expression profiling data paired with CRISPR-induced mutations will be an especially useful source of data of this type, as these experiments include increasingly larger numbers of mutations²². While many of the most useful applications of dimensionality reduction tend to arise from single-cell genomics, for which typical datasets necessitate approaches like UMAP to define relationships between cells, these approaches may also prove useful in visualizing the spatial relationships of biomolecules in tissues²³, genetic interactions, or relationships between human populations²⁴.

Methods

Yeast single-gene deletion transcriptome data. Growth-rate adjusted microarray expression values derived from limma modeling by Kemmeren et al.³, were used as input data. All 1484 single-gene deletion strains from this dataset were used for subsequent dimensionality reduction.

Dimensionality reduction and clustering. To project single-gene deletion strains into two dimensions we performed dimensionality reduction with the UMAP algorithm⁶ using the wrapper function in Monocle 3 (v2.99.3)⁹ to project single-gene deletion strains into two dimensions and subsequently used Louvain clustering²⁵ on strains in 2D UMAP space using default parameters (except, `reduction_method = UMAP`, `metric = cosine`, `n_neighbors = 10`, `min_dist = 0.05`; `clusterCells: method = louvain`, `res = 1e-4`, `k = 3`). Prior to dimensionality reduction, expression values from all 6170 yeast genes were given as input to Principal Component Analysis (PCA) using the Monocle 3 wrapper function “`preprocess_cds`”. The top 100 principal components were then used as input to UMAP for generating 2D projections of the data. For subclustering, main clusters 1–10 were each individually processed using top 25 principal components in the subset data as input to UMAP dimensionality reduction and Louvain clustering (resolution = 1e-4).

Alternative dimensionality reduction with tSNE was performed using the Monocle 3 function `reduceDimension` with default parameters (`reduction_method = tSNE`). Dimensionality reduction using random projection, based on the Johnson-Lindenstrauss lemma, was performed using the `RandPro` (v0.2.0) R package.

Differentially expressed genes per cluster. Gene expression values for single-gene deletions within a cluster were compared to the background set of all deletions. Differentially expressed genes for each cluster were calculated using the `differentialGeneTest()` function in Monocle 3. Because the expression datasets were microarray-derived rather than count-based RNA-seq data, the “`gaussianf`” expression family was used; significance values were corrected for genomic inflation factors using `lambda_g`²⁶.

Benchmarking with known interacting pairs. To test the ability of UMAP distance, and other distance metrics, to capture known interactions, we used a curated consensus set of protein complexes derived from two large, high-throughput mass spectrometry datasets and GO interactions². The consensus set was transformed into a pairwise Boolean interaction matrix based on whether or not each pair had been observed together in the known complex set. Using the subset of pairs that were found in the set of 1484 single-gene deletion transcriptome datasets, for each gene pair, we calculated Euclidean distance in UMAP space. We then used these distances, along with labels for true and false interacting gene pairs derived from gold standard interaction datasets, to generate receiver operating characteristic (ROC) and precision/recall curves with the `PRROC` package in R²⁷.

Reporting summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

The raw data that support the findings of the present study were published previously³ and can be found at <http://deleteome.holstegelab.nl/>. Processed data are available at https://github.com/cole-trapnell-lab/yeast_umap (see Code availability statement).

Code availability

All input data and scripts used for dimensionality reduction and clustering are available through Github (https://github.com/cole-trapnell-lab/yeast_umap).

Received: 7 June 2019; Accepted: 29 February 2020;

Published online: 24 March 2020

References

- Hughes, T. R. et al. Functional discovery via a compendium of expression profiles. *Cell*. [https://doi.org/10.1016/S0092-8674\(00\)00015-5](https://doi.org/10.1016/S0092-8674(00)00015-5) (2000).
- Benschop, J. J. et al. A consensus of core protein complex compositions for *Saccharomyces cerevisiae*. *Mol. Cell* **38**, 916–928 (2010).
- Kemmeren, P. et al. Large-scale genetic perturbations reveal regulatory networks and an abundance of gene-specific repressors. *Cell* **157**, 740–752 (2014).
- Wang, W., Cherry, J. M., Botstein, D. & Li, H. A systematic approach to reconstructing transcription networks in *Saccharomyces cerevisiae*. *Proc. Natl Acad. Sci. USA* **99**, 16893–16898 (2002).
- Laurens van der, Maaten & Hinton, G. Visualizing data using t-SNE. *Laurens. J. Mach. Learn. Res.* **9**, 2579–2605 (2008).
- McInnes, L., Healy, J. & Melville, J. UMAP: Uniform manifold approximation and projection for dimension reduction. *Stat. Mach. Learn. arXiv preprint arXiv:1802.03426* (2018).
- Wang, X. et al. Three-dimensional intact-tissue sequencing of single-cell transcriptional states. *Science*. **361**, eaat5691 (2018).
- Shifrut, E. et al. Genome-wide CRISPR screens in primary human T cells reveal key regulators of immune function resource genome-wide CRISPR Screens in primary human t cells reveal key regulators of immune function. *Cell* **175**, 1958–1971.e15 (2018).
- Cao, J. et al. The single-cell transcriptional landscape of mammalian organogenesis. *Nature* **566**, 496–502 (2019).
- Jean-Baptiste, K. et al. Dynamics of gene expression in single root cells of *A. thaliana*. *Plant Cell*. **31**, 993–1011 (2019).
- Saunders, L. M. et al. Thyroid hormone regulates distinct paths to maturation in pigment cell lineages. *Elife*. <https://doi.org/10.7554/eLife.45181> (2019).
- Guo, L. et al. Resolving cell fate decisions during somatic cell reprogramming by single-cell RNA-Seq. *Mol. Cell* **73**, 815–829 (2019).
- Pijuan-Sala, B. et al. A single-cell molecular map of mouse gastrulation and early organogenesis. *Nature* **566**, 490–495 (2019).
- Becht, E. et al. Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* **37**, 38–47 (2019).
- Schaffrath, R. & Leidel, S. A. Wobble uridine modifications—a reason to live, a reason to die? *RNA Biol.* **14**, 1209–1222 (2017).
- Paraskevopoulou, C., Fairhurst, S. A., Lowe, D. J., Brick, P. & Onesti, S. The Elongator subunit Elp3 contains a Fe4S4 cluster and binds S-adenosylmethionine. *Mol. Microbiol.* **59**, 795–806 (2006).
- Scheidt, V., Juedes, A., Baer, C., Klassen, R. & Schaffrath, R. Loss of wobble uridine modification in tRNA anticodons interferes with TOR pathway signaling. *Microb. Cell* **1**, 416–424 (2014).
- Szklarczyk, D. et al. The STRING database in 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids Res.* **45**, D362–D368 (2017).
- Costanzo, M. et al. A global genetic interaction network maps a wiring diagram of cellular function. *Science*. <https://doi.org/10.1126/science.aaf1420> (2016).
- Cannings, T. I. & Samworth, R. J. Random-projection ensemble classification. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **79**, 959–1035 (2017).
- Bradley, P. H., Gibney, P. A., Botstein, D., Troyanskaya, O. G. & Rabinowitz, J. D. Minor isozymes tailor yeast metabolism to carbon availability. *mSystems* **4**, 1–19 (2019).
- Datlinger, P. et al. Pooled CRISPR screening with single-cell transcriptome readout. *Nat. Methods* **14**, 297–301 (2017).
- Smets, T. et al. Evaluation of distance metrics and spatial autocorrelation in uniform manifold approximation and projection applied to mass spectrometry imaging data. *Anal. Chem.* **91**, 5706–5714 (2019).
- Diaz-Papkovich, A., Anderson-Trocme, L., Ben-Eghan, C. & Gravel, S. UMAP reveals cryptic population structure and phenotype heterogeneity in large genomic cohorts. *PLoS Genet.* **15**, e1008432 (2019).
- Vincent D. Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and E. L. Fast unfolding of communities in large networks. *J. Stat. Mech. Theory Exp.* <https://doi.org/10.1088/1742-5468/2008/10/P10008> (2008).
- Yang, J. et al. Genomic inflation factors under polygenic inheritance. *Eur. J. Hum. Genet.* **19**, 807–812 (2011).
- Grau, J., Grosse, I. & Keilwagen, J. PRROC: computing and visualizing precision-recall and receiver operating characteristic curves in R. *Bioinformatics* **31**, 2595–2597 (2015).

Acknowledgements

We thank J. Packer for advice on differential expression analysis. This work was supported by NIH grants DP2 HD088158, RC2 DK114777, R01HL118342 to C.T.,

GM114166, 1RM1HG010461 to C.Q. and S.F. and P41 GM103533 to S.F. This work was also supported by the Paul G. Allen Frontiers Group.

Author contributions

M.W.D., L.M.S., and C.T. conceived and designed the study. M.W.D. and L.M.S. analyzed the data. M.W.D., L.M.S., and S.F. wrote the paper. C.Q. and C.T. revised the paper.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41467-020-15351-4>.

Correspondence and requests for materials should be addressed to S.F. or C.T.

Peer review information *Nature Communications* thanks the anonymous reviewer(s) for their contribution to the peer review of this work.

Reprints and permission information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020