

Single cell RNA sequencing ecosystem

FASTGenomics: An analytical ecosystem for single-cell RNA sequencing data

Claus J. Scholz¹, Paweł Biernat¹, Matthias Becker¹, Kevin Baßler¹, Patrick Günther¹, Jenny Balfer², Henning Dickten², Lars Flöer², Kathrin Heikamp², Philipp Angerer^{2,3}, Mathias Heilig², Ralf Karle², Meike Köhler², Thomas Mazurkiewicz², Martin Mönninghoff², Christian Sauer², Albrecht Schick², Gerhard Schlemm², Roland Weigelt², Martin Winkler², Thomas Ulas¹, Fabian Theis³, Stephan Huthmacher², Christina Kratsch^{2,*}, Joachim L. Schultze^{1,4,*}

1 Genomics and Immunoregulation, Life & Medical Sciences (LIMES) Institute, University of Bonn, 53115 Bonn, Germany

2 Comma Soft AG, 53229 Bonn, Germany

3 Institute of Computational Biology, German Research Center for Environmental Health, Helmholtz Center Munich, 85764 Munich, Germany

4 Platform for Single Cell Genomics and Epigenomics, German Center for Neurodegenerative Diseases and the University of Bonn, 53175 Bonn, Germany

* Corresponding authors:

Joachim L. Schultze, Genomics and Immunoregulation, Life & Medical Sciences Institute, Carl Troll Strasse 31, 53115 Bonn, Germany, Email: j.schultze@uni-bonn.de, Tel: +49-228-73-62787

Christina Kratsch, Comma Soft AG, Pützchens Chaussee 202, 53229 Bonn, Germany, Email: chistina.kratsch@comma-soft.com

Short title: Single cell RNA sequencing ecosystem

Key words: single cell genomics, single cell RNA-sequencing, ecosystem, systems biology

Single cell RNA sequencing ecosystem

1 **Introduction**

2 Recent technological advances increased the resolution of transcriptomics from cell populations
3 (“bulk”) to single cells¹. While only few cells were assessed in initial projects^{2,3}, evolving technologies
4 now allow the analysis of thousands of cells⁴⁻⁶, with the largest publicly available dataset currently
5 comprising more than 1.3 million cells⁷. In contrast to bulk RNA sequencing (RNA-seq), single cell (sc)
6 technologies are much more demanding due to high technical variation with zero-inflation being a
7 major property⁸. As a consequence, a myriad of novel computational approaches and tools have been
8 developed for the different scRNA-seq technologies⁹, but these thriving innovations also constitute a
9 lack of widely accepted gold standards for data analysis. By construction, many of the proposed
10 algorithms and approaches address only certain steps in the analytical scRNA-seq workflow, are
11 adapted to certain scRNA-seq technologies, or cannot be easily combined with other tools, limiting
12 their broad applicability. Notable exceptions are software packages like Monocle¹⁰, Seurat¹¹ and
13 Scanpy¹², which are well documented, cover big parts of the analysis workflow, and are flexible in their
14 application; nevertheless, due to their command line-based environments, they are still restricting
15 access to scRNA-seq for the broader life and medical sciences community. More user-friendly tools
16 with graphical user interfaces have been introduced, like Granatum¹³, which offers a local installation,
17 or the online tool ASAP¹⁴ and the commercial solution SeqGeq¹⁵. In their current versions, they offer
18 popular analysis algorithms, yet are limited in scalability in a multi-user setting, data security, usability
19 of data varying in size over several orders of magnitude, and integration of own analytical concepts.
20 Especially for largest-scale single cell genomics undertakings like the Human Cell Atlas (HCA)¹⁶, existing
21 tools provide only limited analytical performance due to inefficient resource allocation for exploding
22 memory and computing requirements for datasets in the magnitude of millions of cells, thus
23 underscoring the necessity for a powerful software solution tailored to efficiently handle mega-
24 analyses through distributed computing.

25 Single cell genomics - with scRNA-seq leading the way - will revolutionize the life and medical
26 sciences^{8,17-19}. Here, we postulate that an analytical ecosystem for single cell genomics applications
27 will foster research and development in this field. Such an ecosystem should give computational
28 experts a platform to make their tools available to a broader audience in a user-friendly fashion, allow
29 high-end users to develop individualized workflows, and provide the novice user a computational
30 environment to get acquainted with the special computational requirements for single cell analysis.
31 Furthermore, such an ecosystem should serve as a platform for the community to share public datasets
32 with a broader audience by following the FAIR Guiding Principles for scientific data management and
33 stewardship²⁰, provide a scalable infrastructure for projects with large datasets even across numerous
34 institutions, host benchmarking capabilities for newly developed algorithms for the analysis of scRNA-
35 seq data, and even serve as a portal for large international projects such as the HCA¹⁶. Finally, an
36 analytical ecosystem must implement best practice measures that agree with institutional and
37 governmental data security regulations. To address all these requirements, we have developed
38 FASTGenomics (<https://fastgenomics.org>) as a powerful, efficient, versatile, robust, safe and intuitive
39 analytical ecosystem for single-cell transcriptomics. Access to the FASTGenomics ecosystem and its
40 functionality is granted for free upon registration to allow unrestricted interaction with the single cell
41 genomics community and especially academia. Furthermore, as suggested by the HCA white paper and
42 guided by representatives of the HCA, the implementation of FASTGenomics as a portal for the HCA is
43 currently on its way.

44
45

Single cell RNA sequencing ecosystem

46 **App Store in FASTGenomics serves as platform for novel algorithms**

47 At the heart of FASTGenomics is a hybrid app store (**Figure 1A**) optionally composed of public (cloud)
48 and private (local) app repositories hosting algorithms for calculations and data visualization. Novel
49 algorithms can be provided as new apps by the computational biology community (**Figure 1B**). The
50 well-documented application program interface (API) (**Supplementary Information “Description of
51 the API of FASTGenomics”**) defines data input and output (**Figure 1C**) and allows seamless integration
52 into the FASTGenomics ecosystem. Apps submitted to the public app repository
53 (<https://github.com/fastgenomics>) are included in the complete end-user environment
54 (**Supplementary Information “Detailed description of end-user experience of the FASTGenomics
55 ecosystem”**). Additionally, designing customized workflows integrating custom-made apps is a major
56 feature of FASTGenomics (**Figure 1D**). Furthermore, the workflow editor allows to adjust
57 parametrization of apps, thus providing a maximum of analytical flexibility. Currently, workflow editing
58 is done via the command line, the next version of the workflow editor is planned to provide an intuitive
59 graphical user interface with functionality to share custom workflows (**Supplementary Figure S1**).

60

61 **Architecture, scalability and data security of the Docker-based hybrid model of FASTGenomics**

62 The FASTGenomics ecosystem has been implemented as a Docker²¹-based cloud solution, which can
63 also be used as a local environment with a community-wide app repository (hybrid design) allowing to
64 share data, apps and workflows, but also information, expertise and knowledge about single cell
65 genomic analyses (**Figure 1A**, for a user perspective, **Supplementary Figure S2** for architectural
66 specifications, for more details see **Supplementary Information “Technical realization of
67 FASTGenomics with Docker-based cloud solution”**). Alternatively, entirely local installations – as they
68 might be required within industry – are also possible. While ensuring standardization and reduced
69 administrative burden, the modular, docker-based hybrid cloud solution of FASTGenomics also
70 provides the necessary scalability to run projects with very large datasets. A dynamic allocation and
71 flexible use of available resources will be achieved by leveraging Kubernetes technology in the next
72 release of the platform^{22,23}.

73 In its current version, FASTGenomics is being developed according to EU-GDPR (General Data
74 Protection Regulation) and the German Federal Data Protection Act (“Bundesdatenschutzgesetz”,
75 BDSG), one of the strictest data protection laws in the world. To minimize security issues related to
76 multi-user access to the platform and the use of custom apps, FASTGenomics implements a rigorous
77 multi-layer security concept of data encryption, controlled access and transfer to protect study data
78 (expression tables, sample metadata and analysis results) as well as user data from unauthorized
79 access and manipulation (**Supplementary Figure S3**). A data protection concept has been developed
80 accordingly and will be continuously updated according to legal requirements (**Supplementary
81 Information “Data Security Concept within FASTGenomics”**).

82

83 **User-friendly computational environment**

84 Within the FASTGenomics ecosystem, analyses can be initiated and monitored from essentially any
85 web-compatible hardware with a web browser, without requiring extensive computing or memory
86 resources locally. For the end-user following registration, FASTGenomics provides an interface for data
87 upload (**Supplementary Figure S4A**, **Supplementary Information “Description of data upload via
88 upload Dock in FASTGenomics”**), starting from count tables and experimental metadata, followed by
89 standardized quality checks, e.g. average molecule counts, gene types, and quantification of batch

Single cell RNA sequencing ecosystem

90 effects (**Supplementary Figure S4C-E**), and two pre-defined data analysis and visualization workflows,
91 ‘Subtype Discovery’ and ‘Pseudo Time Analysis’ (**Figure 1D**). The former includes a neural network
92 approximation of the parametric tSNE²⁴ and a 3D visualization of cells with coloring according to cluster
93 assignments, gene expression and metadata (see also **Supplementary Table 1**). Each analysis results
94 in the definition of genes of interest and a functional categorization with the help of external
95 databases, e.g. Gene Ontology (GO²⁵). Workflows in FASTGenomics end with a summary, a detailed
96 description of all analysis steps including information about algorithms, software, versions, and
97 parametrizations used as well as input data and results produced (**Figure 1D, Supplementary Figure**
98 **S5A, S5B, Supplementary Information “Description of Summary of any given analysis”**). The summary
99 is intended to maximize reproducibility and transparency of the analysis, which could be made
100 available e.g. in scientific publications or within documentation required in regulatory environments.

101

102 **Platform for sharing datasets for further public exploitation**

103 Another important feature of FASTGenomics is a standardized package for public dataset presentation,
104 which we utilized to present 10 recently published datasets ranging from 482 to 68,579 cells per
105 dataset (**Supplementary Table 2**)^{5,26–34}. Available datasets can be connected with standard workflows
106 provided by FASTGenomics, but also with customized apps and workflows as exemplified for a previous
107 MARS-Seq dataset (**Supplementary Figure S6**)³¹. By combining a dataset with the initial analysis, the
108 data can be examined by anybody following the same algorithmic settings as previously reported in
109 the literature. Moreover, this also allows to compare different analysis strategies directly on the same
110 platform. We also performed concordance analyses for selected datasets presented in FASTGenomics
111 (**Figure 2A**) and focus here on a dataset with 3,005 cells published by Zeisel et al.³⁴. Using the BACKSPIN
112 clustering algorithm, a total of 9 clusters that were assigned to 7 classes of cell types were previously
113 identified in the dataset, while after our neural network-based dimensionality reduction a subset of
114 2,375 cells could be assigned to 16 clusters. Thus, the FASTGenomics ‘subtype discovery’ standard
115 workflow revealed a more fine-grained cluster structure than the BACKSPIN algorithm while preserving
116 the co-clustering of functionally closely related cell types. In particular, neuronal and glial cell types
117 were clearly distinguished from each other as well as from vasculature; in more detail,
118 oligodendrocytes and pyramidal neurons were each assigned to one FASTGenomics cluster, while
119 interneurons were clustered to six main classes. Quantitatively this translates to an adjusted mutual
120 information value of 0.75 and median concordance rates of 96.5% for FASTGenomics and 90% for
121 BACKSPIN (**Figure 2AB, Supplementary Information**). Such measures might be also used to estimate
122 specialized analyses settings in previously published datasets. Collectively, the option to freely share
123 previously published large datasets on FASTGenomics allows intuitive and interactive cross-
124 examination, which goes far beyond the current options in scientific publications.

125

126 **FASTGenomics provides higher flexibility and scalability compared to existing platforms**

127 Next, we intended to compare FASTGenomics to the three currently available GUI-based platforms
128 ASAP¹⁴, Granatum¹³ and SeqGeq¹⁵ (for detailed setup see **Supplementary Information “Setup of ASAP,**
129 **Granatum and SeqGeq for comparison with FASTGenomics”**). We utilized five datasets ranging from
130 1,920³³ to 68,579 cells²⁹ and compared for data upload, pre-processing cell clustering, differential gene
131 expression analysis, pseudo time analysis and analysis summary. In their default configuration, among
132 the four evaluated tools, only FASTGenomics performed all steps with all datasets (**Figure 2C**). We
133 furthermore determined the resources needed by FASTGenomics to compute analyses with different

Single cell RNA sequencing ecosystem

134 dataset sizes (experiment details in “**Resource Requirements of a FASTGenomics Analysis**
135 **Workflow**”). Analysis runtime and memory requirements are both strongly correlated and depend on
136 the number of cells analyzed; furthermore, analysis of all datasets across the tested size range is
137 feasible with a contemporary desktop computer (**Figure 2D**).

138

139 **Outlook**

140 In upcoming versions of FASTGenomics, datasets, apps and workflows can be shared in private
141 spaces/sections between collaboration partners prior to publishing, thus providing the infrastructure
142 for multi-institutional collaboration projects. Furthermore, import/export apps will be implemented
143 to be fully interoperable with established analysis software tools like Monocle¹⁰, Scanpy¹², Scater³⁵,
144 Seurat¹¹, etc., but also with data repositories like Gene Expression Omnibus (GEO)³⁶. Finally, a
145 connection of FASTGenomics to major laboratory information management systems (LIMS) for the
146 import of experimental variables as metadata for new datasets as well as the export of the analysis
147 summary back to the experimenters’ LIMS is currently evaluated and discussed with future users.

148

149 **Conclusion**

150 Taken together, FASTGenomics is designed as a secure, flexible, scalable but also standardized
151 platform for single cell RNA-seq data, open to the scientific community. A major feature is to provide
152 highest reproducibility and transparency for single cell data analysis to the whole community. Due to
153 its modular and open structure it could also serve as a platform for community-wide benchmarking for
154 novel algorithms and even serve as one of the tertiary portals planned within the HCA data
155 coordination platform of the Human Cell Atlas¹⁶. Furthermore, by design, it scales already routinely to
156 more than 5×10^4 cells per project and prototype apps suggest that scaling to 10^6 cells is also possible.
157 Moreover, its hybrid design will also allow using FASTGenomics on premise, which might be of interest
158 to clinical research and the pharmaceutical industry.

159

160 **Figure Legends**

161 **Figure 1: FAST Genomics ecosystem. (A)** Hybrid app store concept. To provide both the advantages of
162 community access to the FASTGenomics framework as well as the security of a private working
163 environment, FASTGenomics runs in the cloud and can also be installed on premise. The cloud
164 installation allows the usage of public apps and exchange with the global research community, whereas
165 the on-premise installation could run on a local cluster. Additional local app repositories and data
166 storage can be added for private access only. **(B)** Typical structure of a FASTGenomics workflow. All
167 FASTGenomics workflows consist of calculation apps (such as quality checks, data normalization,
168 dimensionality reduction, clustering, ...) that take inputs and consecutively produce new results for
169 upstream calculation apps. Selected outputs of the calculation workflow are displayed in the browser
170 with the help of visualization apps in the according visualization workflow. **(C)** Structure of a
171 FASTGenomics app. Apps are Docker containers that interact with the FASTGenomics framework using
172 an interface for data input and a configuration file providing necessary parameters for the analysis.
173 Each FASTGenomics app dynamically generates a summary of the analysis performed by the app that
174 is collected by the FASTGenomics summary service. Depending on app type, different channels are
175 used for results, calculation apps write output to disk, whereas visualization apps send output to the

Single cell RNA sequencing ecosystem

176 web browser. The use of the Docker framework enables app developers to implement algorithms in
177 any programming language of choice. A detailed tutorial for the development of calculation and
178 visualization apps as well as sample code can be found at the public FASTGenomics app repository
179 (<https://github.com/fastgenomics>). **(D)** Workflow definitions and app concept: workflow definitions
180 are configuration files that describe the calculation and visualization apps used for a specific workflow.
181 User-defined workflows can be added simply by creating new workflow definitions, which may recycle
182 previously defined apps. In particular, apps for the exploration of gene candidate lists with the help of
183 DE analysis and functional annotation are typical candidates for multi-workflow apps. All workflows
184 end with a detailed summary of the analyses performed to ensure maximum transparency and
185 reproducibility.

186

187 **Fig 2: Reproducibility of workflows and performance of FASTGenomics. (A)** Clustering results of
188 individual cells generated with the standard 'subtype discovery' workflow in FASTGenomics were
189 compared to published findings by determination of the adjusted mutual information (AMI). Immune
190 cell datasets^{29,31} displayed a lower degree of concordance than neuronal³⁴, cancer²⁷ and retinal tissue⁵
191 datasets, presumably due to the lower RNA content of immune cells²⁹ and the lower number of genes
192 expressed²⁷. **(B)** FASTGenomics (FG) cluster assignments compared to published cell types³⁴. For each
193 FG cluster, the proportion of main cell types (inner circle) and subtypes (outer circle) are shown. The
194 FASTGenomics standard 'subtype discovery' workflow clearly distinguished single-cell transcriptomes
195 at higher resolution than main cell types, but with lower resolution than the published subclustering
196 approach. Based on single-cell transcriptomic data, biologically meaningful subclasses were generated
197 by the FASTGenomics 'subtype discovery' workflow, classifying neuronal and glial cells, vasculature
198 and immune cell types in distinct units. **(C)** Performance comparison between FASTGenomics and three
199 additional GUI-based platforms for single cell analysis. FASTGenomics (<https://fastgenomics.org>) was
200 compared to the online tool ASAP (<https://asap.epfl.ch/>) and local installations of Granatum
201 (<http://garmiregroup.org/granatum/app>) and SeqGeq (<https://www.flowjo.com/solutions/seqgeq>)
202 installed on a 64 bit Windows 10 machine with Intel i7 6700K CPU and 32 GB RAM). Comparison was
203 performed in 7 categories (data upload, data preprocessing, cell clustering, differential gene
204 expression, functional analysis, pseudotime analysis, analysis summary). Datasets of various sizes,
205 ranging from 1,920 to 68,579 cells^{5,29,32-34} were used to assess scalability of the platforms. The size of
206 the largest dataset, for which an analysis task could be accomplished is shown for all evaluated
207 pipelines. **(D)** Required resources for analysis of data sets of various sizes^{5,29,32-364}. Maximum memory
208 usage (blue dots) and overall analysis runtime (red dots) to complete data normalization,
209 dimensionality reduction and cell clustering are shown depending on the number of cells contained in
210 each analyzed dataset.

211

Single cell RNA sequencing ecosystem

212 **References**

- 213 1. Picelli, S. Single-cell RNA-sequencing: The future of genome biology is now. *RNA Biol.* 1–14
214 (2016). doi:10.1080/15476286.2016.1201618
- 215 2. Shalek, A. K. *et al.* Single-cell transcriptomics reveals bimodality in expression and splicing in
216 immune cells. *Nature* **498**, 236–240 (2013).
- 217 3. Islam, S. *et al.* Characterization of the single-cell transcriptional landscape by highly multiplex
218 RNA-seq. *Genome Res.* **21**, 1160–1167 (2011).
- 219 4. Jaitin, D. A. *et al.* Massively parallel single-cell RNA-seq for marker-free decomposition of tissues
220 into cell types. *Science* **343**, 776–779 (2014).
- 221 5. Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual Cells Using
222 Nanoliter Droplets. *Cell* **161**, 1202–1214 (2015).
- 223 6. Gierahn, T. M. *et al.* Seq-Well: portable, low-cost RNA sequencing of single cells at high
224 throughput. *Nat. Methods* **14**, 395–398 (2017).
- 225 7. 10X Genomics. Megacell 1.3 Mio dataset. (2017). Available at:
226 [https://community.10xgenomics.com/t5/10x-Blog/Our-1-3-million-single-cell-dataset-is-ready-](https://community.10xgenomics.com/t5/10x-Blog/Our-1-3-million-single-cell-dataset-is-ready-to-download/ba-p/276)
227 [to-download/ba-p/276](https://community.10xgenomics.com/t5/10x-Blog/Our-1-3-million-single-cell-dataset-is-ready-to-download/ba-p/276). (Accessed: 6th November 2017)
- 228 8. Stegle, O., Teichmann, S. A. & Marioni, J. C. Computational and analytical challenges in single-cell
229 transcriptomics. *Nat. Rev. Genet.* **16**, 133–145 (2015).
- 230 9. Poirion, O. B., Zhu, X., Ching, T. & Garmire, L. Single-Cell Transcriptomics Bioinformatics and
231 Computational Challenges. *Front. Genet.* **7**, 163 (2016).
- 232 10. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
233 pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
- 234 11. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell
235 gene expression data. *Nat. Biotechnol.* **33**, 495–502 (2015).
- 236 12. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data
237 analysis. *Genome Biol.* **19**, 15 (2018).

Single cell RNA sequencing ecosystem

- 238 13. Zhu, X. *et al.* Granatum: a graphical single-cell RNA-Seq analysis pipeline for genomics scientists.
239 *Genome Med.* **9**, 108 (2017).
- 240 14. Gardeux, V., David, F. P. A., Shajkofci, A., Schwalie, P. C. & Deplancke, B. ASAP: a Web-based
241 platform for the analysis and interactive visualization of single-cell RNA-seq data. *Bioinforma.*
242 *Oxf. Engl.* (2017). doi:10.1093/bioinformatics/btx337
- 243 15. SeqGeq® | FlowJo, LLC. Available at: <https://www.flowjo.com/solutions/seqgeq>. (Accessed: 2nd
244 February 2018)
- 245 16. Human Cell Atlas. The Human Cell Atlas White Paper. Available at:
246 <https://www.humancellatlas.org/news/13>. (Accessed: 6th November 2017)
- 247 17. Shapiro, E., Biezuner, T. & Linnarsson, S. Single-cell sequencing-based technologies will
248 revolutionize whole-organism science. *Nat. Rev. Genet.* **14**, 618–630 (2013).
- 249 18. Jaitin, D. A., Keren-Shaul, H., Elefant, N. & Amit, I. Each cell counts: Hematopoiesis and immunity
250 research in the era of single cell genomics. *Semin. Immunol.* **27**, 67–71 (2015).
- 251 19. Wagner, A., Regev, A. & Yosef, N. Revealing the vectors of cellular identity with single-cell
252 genomics. *Nat. Biotechnol.* **34**, 1145–1160 (2016).
- 253 20. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
254 stewardship. *Sci. Data* **3**, 160018 (2016).
- 255 21. Beaulieu-Jones, B. K. & Greene, C. S. Reproducibility of computational workflows is automated
256 using continuous analysis. *Nat. Biotechnol.* **35**, 342–346 (2017).
- 257 22. Burns, B., Grant, B., Oppenheimer, D., Brewer, E. & Wilkes, J. Borg, Omega, and Kubernetes.
258 *Queue* **14**, 10:70–10:93 (2016).
- 259 23. Schulz, W. L., Durant, T. J. S., Siddon, A. J. & Torres, R. Use of application containers and
260 workflows for genomic data analysis. *J. Pathol. Inform.* **7**, (2016).
- 261 24. Maaten, L. Learning a Parametric Embedding by Preserving Local Structure. in *PMLR* 384–391
262 (2009).
- 263 25. Gene Ontology Consortium. Gene Ontology Consortium: going forward. *Nucleic Acids Res.* **43**,
264 D1049-1056 (2015).

Single cell RNA sequencing ecosystem

- 265 26. Tirosh, I. *et al.* Single-cell RNA-seq supports a developmental hierarchy in human
266 oligodendroglioma. *Nature* **539**, 309–313 (2016).
- 267 27. Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-
268 seq. *Science* **352**, 189–196 (2016).
- 269 28. Ziegenhain, C. *et al.* Comparative Analysis of Single-Cell RNA Sequencing Methods. *Mol. Cell* **65**,
270 631–643.e4 (2017).
- 271 29. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat.*
272 *Commun.* **8**, 14049 (2017).
- 273 30. Moignard, V. *et al.* Decoding the regulatory network of early blood development from single-cell
274 gene expression measurements. *Nat. Biotechnol.* **33**, 269–276 (2015).
- 275 31. Mass, E. *et al.* Specification of tissue-resident macrophages during organogenesis. *Science* **353**,
276 (2016).
- 277 32. Paul, F. *et al.* Transcriptional Heterogeneity and Lineage Commitment in Myeloid Progenitors.
278 *Cell* **163**, 1663–1677 (2015).
- 279 33. Nestorowa, S. *et al.* A single-cell resolution map of mouse hematopoietic stem and progenitor
280 cell differentiation. *Blood* **128**, e20-31 (2016).
- 281 34. Zeisel, A. *et al.* Brain structure. Cell types in the mouse cortex and hippocampus revealed by
282 single-cell RNA-seq. *Science* **347**, 1138–1142 (2015).
- 283 35. McCarthy, D. J., Campbell, K. R., Lun, A. T. L. & Wills, Q. F. Scater: pre-processing, quality control,
284 normalization and visualization of single-cell RNA-seq data in R. *Bioinforma. Oxf. Engl.* **33**, 1179–
285 1186 (2017).
- 286 36. Barrett, T. *et al.* NCBI GEO: archive for functional genomics data sets--update. *Nucleic Acids Res.*
287 **41**, D991-995 (2013).
- 288
- 289

Single cell RNA sequencing ecosystem

290 **Acknowledgements**

291 We would like to thank John Marioni for fruitful interactions concerning the project to implement
292 FASTGenomics as one of the portals for the HCA. This work was supported by a grant from the Federal
293 Ministry for Economic Affairs and Energy (BMW i Project FASTGENOMICS). JLS is a member of the
294 excellence cluster immunosensation. JLS is further supported by the DFG (Sachbeihilfe SCHU 950/9-1;
295 SFB 704, projects A13, Z5; excellence cluster immunosensation).

296

297 **Author contributions**

298 CJS wrote the manuscript; CJS, PB, MB, KB, PG, JB, HD, LF, KH, PA, MK and TU designed analyses; CJS,
299 PB, JB, HD, LF, KH, PA, MK, MM, CS, AS and GS implemented apps; HD, MH, RK, TM, MM, CS, AS, GS,
300 RW and MW developed the platform; JB, HD, KH, RK, TM and CS contributed to the manuscript; KH
301 and JB managed the project; SH perceived idea, managed and supervised the project; CK managed and
302 supervised project, wrote the manuscript; JLS, FT discussed and improved the project and the
303 manuscript; JLS perceived idea for the project, managed and supervised the project, wrote and
304 designed the manuscript.

305

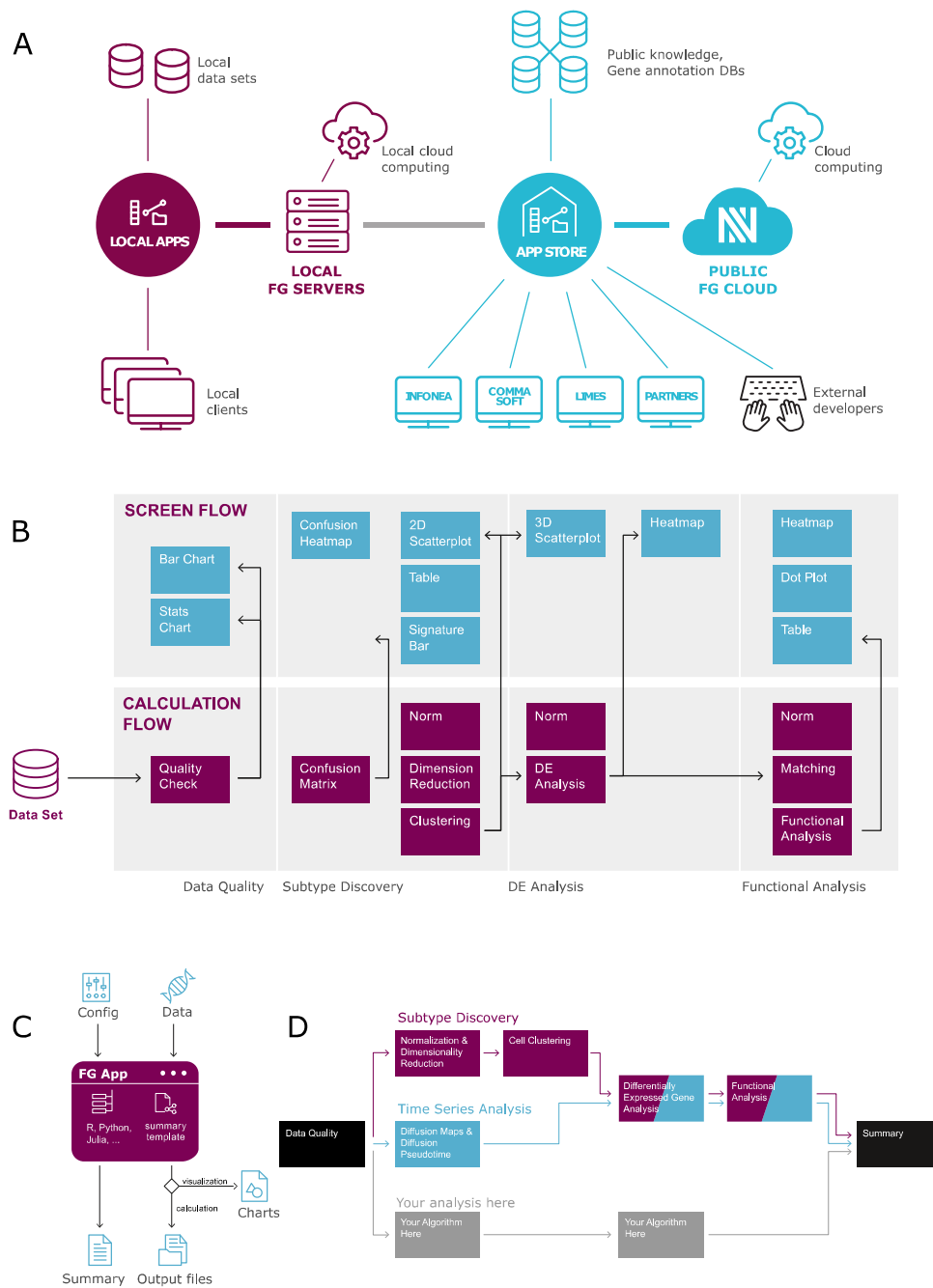
306 **Competing financial interests**

307 CJS, PB, MB, KB, PG, TU, FT and JLS declare no competing financial interests. PA, JB, HD, LF, KH, MH,
308 RK, MK, TM, MM, CS, AS, GS, MW, SH and CK are paid employees of Comma Soft AG, a commercial
309 company developing the FASTGenomics platform.

310

Single cell RNA sequencing ecosystem

Figure 1



311

312

Single cell RNA sequencing ecosystem

Figure 2

