

1 **SUPPLEMENTARY INFORMATION, FIGURES AND TABLES**

2

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

4

5

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

11

12 1 Genomics and Immunoregulation, Life & Medical Sciences (LIMES) Institute, University of Bonn, 53115
13 Bonn, Germany

14 2 Comma Soft AG, 53229 Bonn, Germany

15 3 Institute of Computational Biology, German Research Center for Environmental Health, Helmholtz
16 Center Munich, 85764 Munich, Germany

17 4 Platform for Single Cell Genomics and Epigenomics, German Center for Neurodegenerative Diseases
18 and the University of Bonn, 53175 Bonn, Germany

19 * Corresponding authors:

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

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

24

25

26 Short title: Single-cell RNA sequencing ecosystem

27

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

29

30 **Supplementary Information**

31 **Content:**

- 32 • **Description of the API of FASTGenomics**
- 33 • **Detailed description of end-user experience of the FASTGenomics ecosystem**
- 34 • **Technical realization of FASTGenomics with Docker-based cloud solution**
- 35 • **Logical implementation built on the Docker-based cloud solution**
- 36 • **Adding third party apps to the Docker-based cloud solution**
- 37 • **Data Security Concept within FASTGenomics**
- 38 • **Description of data upload to FASTGenomics**
- 39 • **Description of Summary of any given analysis**
- 40 • **Concordance rate analysis of FASTGenomics Analyses with published results**
- 41 • **Neural network-based dimensionality reduction and clustering**
- 42 • **Methods applied to determine concordance**
- 43 • **Results of concordance rate analysis**
- 44 • **Setup of ASAP, Granatum and SeqGeq for comparison with FASTGenomics**
- 45 • **Supplementary Figures 1 to 6 (including legends)**
- 46 • **Supplementary Tables 1 and 2**
- 47

48 Description of the API of FASTGenomics

49 The FASTGenomics ecosystem for single-cell analyses allows integrating algorithms from third parties. For
50 proper integration into the FASTGenomics pipeline, implementations must comply with the
51 FASTGenomics application-programming interface (API). This requires that apps are provided as Docker
52 containers of defined structure, with version information provided in Docker tags. The root directory must
53 include a `Dockerfile`, the script to prepare and invoke app components. A `manifest.json` file that
54 contains app descriptions, input, output and parameter definitions. The `sample_data` folder containing
55 files required for integrity tests during app validation. The `readme.md` file with the app documentation
56 and the source code providing the app functionality. Further optional or best practice components of a
57 generic FASTGenomics app include the `docker-compose.yml` file with information how to build and
58 start the Docker container and providing input/output directories. A `requirements.txt` detailing
59 dependencies for proper app functioning. Markdown-formatted template file(s) for the summary output
60 that is dynamically filled with information generated during the app run. FASTGenomics distinguishes
61 between two types of apps that interact with different parts of the pipeline, calculation apps are handled
62 by the workflow engine, while visualization apps are managed by the workflow client. Both app types use
63 defined mountpoints for data (`/fastgenomics/data`, read-only) and analysis configuration file input
64 (`/fastgenomics/config`, read-only). Calculation apps write results (`/fastgenomics/output`,
65 read/write) and a summary (`/fastgenomics/summary`, read/write) to disk, whereas visualization
66 apps send output to the web browser via port 8000. Apps developed to comply with the FASTGenomics
67 API can be published in the public FASTGenomics app repository (<https://github.com/fastgenomics>) and
68 will be made available in the FASTGenomics Docker registry. A detailed tutorial for the development of
69 calculation and visualization apps as well as sample code can be found at
70 <https://github.com/fastgenomics>.

71

72 Detailed description of end-user experience of the FASTGenomics ecosystem

73 FASTGenomics (<https://fastgenomics.org>) allows several distinct levels of usage and access. An
74 anonymous access allows users interested in FASTGenomics to see results from pre-calculated analyses
75 on a selection of publicly available datasets. Full access to all functionality requires a free user registration,
76 which enables usage of a greater set of experiments from public repositories, the possibility to integrate
77 own data sets, as well as the availability of the full range of analysis tools. The natural first step of a new
78 project is to upload single-cell expression data to FASTGenomics after which the data may be checked for
79 quality by exclusion of cells with too few expressed genes, and exclusion of genes expressed in too few
80 cells (**Supplementary Figure S4E**). The according calculation app can be incorporated into the workflows
81 to perform this task (**Supplementary Table 1A**). This step is followed by the data quality screen showing
82 general statistics about the data, such as average molecule counts and quantification of batch effects. For
83 data analysis, FASTGenomics offers two pre-defined alternatives: subtype discovery and time series
84 analysis (**Figure 1D**). Furthermore, the workflow editor allows creating custom analytical scenarios that
85 can involve any app available in the FASTGenomics app store. In both pre-defined workflows, an overview
86 of the dataset detailing aspects of data quality (e.g. summary statistics on expression values, presence of

87 putative batch effects, etc.) is given in the first screen (**Supplementary Figure S4B-D**). The subtype
88 discovery workflow proceeds with data normalization and dimensionality reduction, followed by
89 clustering of cells (**Supplementary Figure S1C,D**). This cluster projection is displayed in an aquarium plot,
90 with the first two dimensions corresponding to coordinates determined in our parametric t-SNE approach
91 and the third dimension representing the cluster assignment confidence (with high-confidence cluster
92 assignments “swimming” on top and low-confidence assigned cells sinking to the ground). The next step
93 detects differentially expressed genes between clusters and display these in a heatmap (data not shown).
94 In the pseudotime workflow, diffusion maps are generated to order cellular transcriptomes along
95 pseudotemporal axes. This workflow also determines genes responsible for branches in the trajectory.
96 Both, the subtype discovery as well as the pseudotime workflow conclude the analytic sequence with the
97 functional characterization of signature genes using gene ontology enrichments. At the end of all
98 workflows, a detailed summary is dynamically generated during runtime and displayed to the user
99 (**Supplementary Figure S5B**).

100

101 **Technical realization of FASTGenomics with Docker-based cloud solution**

102 The publicly accessible instance of the FASTGenomics ecosystem (<https://fastgenomics.org>) is installed in
103 the Microsoft Azure cloud hosted by server infrastructure located in Western Europe. Currently,
104 FASTGenomics services run on a Standard D8s v3 system (8 vCPUs based on the 2.3 GHz Intel XEON[®] E5-
105 2673 v4 processor and 32 GB RAM). However, other options can be envisioned since the system is
106 designed to allow hybrid computing by integrating local and public cloud installations.

107 The FASTGenomics ecosystem itself is based on Docker infrastructure (currently using version 17.06.0-ce,
108 build 02c1d87), all pipeline components as well as calculation and visualization apps are packaged in
109 Docker containers. Internal pipeline components and apps are deployed to different Docker registries, the
110 former can only be accessed by the runtime environment and FASTGenomics administrators, while the
111 latter is also accessible to the community to allow contribution of apps.

112 The FASTGenomics ecosystem consists of several major components, each with its own responsibility as
113 shown in **Supplementary Figure 2**. The user directly interacts with the `FASTGenomics Client` via
114 the web browser. This serves to display the FASTGenomics website, where users can login, upload and
115 select datasets, choose analysis workflows and access other specialized services. A calculation engine
116 consisting of the `Workflow Engine`, the `Task Dispatcher` and the `Container Service`
117 manages the application of an analysis workflow on a selected dataset. The `Workflow Engine` uses
118 the `Container Service` to start and stop calculation apps that perform individual analysis steps. The
119 latter reports the calculation status to the `Task Dispatcher`, that initializes (creates a unique ID) and
120 finalizes analysis instances upon requests by the `Workflow Engine`. The `Workflow Client` shows
121 a screen flow in the web browser to display results from analyzes; this also makes use of the `Container`
122 `Service` to start and stop visualization apps. Both components access the `Data Store`, which
123 organizes data management for an analysis. Finally, the data upload system consists of the `Upload`
124 `Client`, which takes care of data transfer from the user’s system to the FASTGenomics servers, as well

125 as the `Packaging Service` integrating the uploaded data into the FASTGenomics system. Overall,
126 connection to the user is secured by an `OpenID Connect` component to ensure that only validated
127 users can access the application.

128

129 **Logical implementation built on the Docker-based cloud solution**

130 Each FASTGenomics project consists of one or more analyses and the data set that is to be analyzed. An
131 analysis in turn can be broken down into the algorithms that calculate results and the visualizations that
132 display these results (**Supplementary Figure S2**). In FASTGenomics, the former are called calculation apps
133 and are combined into workflows, while the latter are called visualization apps and encapsulated into
134 screen flows.

135

136 **Adding third party apps to the Docker-based cloud solution**

137 The scientific community may develop individual apps for FASTGenomics, which are wrapped up in Docker
138 images. A provider of apps is authorized to push Docker images to the public FASTGenomics Docker
139 registry where the images are retrieved on demand by the system. Apps can be calculations producing
140 results or visualizations displaying results. To use a custom app in a FASTGenomics analysis, workflow or
141 screen flow definitions are adjusted. In the current online release, two pre-defined analysis workflows are
142 integrated. Adding custom apps and analyses is a feature of FASTGenomics that is predicted to be active
143 and under continuous development.

144

145 **Data Security Concept within FASTGenomics**

146 The FASTGenomics ecosystem as well as any other web-accessible multi-user platform storing and
147 analyzing sensitive data (e.g. unpublished experimental data, clinically relevant metadata, user data) are
148 subject to tight regulations for data security. As maintainer of an online platform, FASTGenomics needs
149 to adhere to the law including the German Federal Data Protection Act (“Bundesdatenschutzgesetz”,
150 BDSG¹) and of May 2018 the European General Data Protection Regulation (GDPR) (Regulation (EU)
151 2016/679²). These regulations cover diverse aspects of data management and safety.

152 In order to respond to these regulations, FASTGenomics is continuously working on a security concept
153 defining the essential, recommended and desirable security features of a single cell analysis platform. This
154 is an ongoing process to account for new developments and planned future components and
155 functionality. At the current state, FASTGenomics has the following security features implemented:

- 156 • All data in FASTGenomics are stored on encrypted data volumes.
- 157 • To ensure safe network topology, both the external and the internal communication between the
158 components is encrypted by HTTPS.

- 159
- 160
- 161
- 162
- 163
- 164
- 165
- 166
- 167
- 168
- 169
- 170
- 171
- 172
- 173
- Authentication is achieved by an OpenID Connect Provider. After registration, users are required to confirm their identity via mail to avoid platform misuse by bots.
 - While accessing the platform and user data therein, access is again regulated via authorization checking done by each involved application. This feature manages access rights of each user, for example when retrieving information from the data module. Here, the platform ensures that only eligible private and public data sets are visible for the current user (**Supplementary Figure S3**).
 - Finally, several security features address the setup of the Docker container making up the FASTGenomics infrastructure. No Docker container is allowed to have root access. Containers that communicate with external components enforce authenticated users and only communicate using HTTPS. The export of a port is limited to this container group.
 - Implementation of national legal requirements for intellectual property with respect to software development (app development).
 - Definition and monitoring of organizational best practices for all processes involving data handling, resource access allocation and platform administration.
 - Definition and implementation of best practices for internet access of apps.

174 In addition to these security features, the FASTGenomics security concept addresses further actions for
175 risk minimization and data protection as features planned for future development. These include:

- 176
- 177
- 178
- 179
- 180
- 181
- 182
- 183
- 184
- 185
- implementation of a framework for error logging, data access, and data manipulation
 - rules that provide manipulation security of data and apps
 - definition of best practices in the context of software development in general and the use of container frameworks like Docker in particular, e.g. managing resource access of apps
 - definition and implementation of rules for computing resource access of apps
 - software support for complete data removal upon user request
 - definition and implementation of best practices for validation of usage statistics and application of web tracking software (e.g. Google Analytics)
 - definition and implementation of best practices for anonymous access to suitable resources
 - definition and implementation of best practices for data publication

186

187 Apart from ensuring data security, FASTGenomics aims to provide a good framework for the
188 reproducibility of analyses and the sharing of data and knowledge. Such aspects are increasingly discussed
189 in the scientific world, and driven by concepts like FAIR aiming to facilitate research and knowledge
190 discovery by Findable, Accessible, Interoperable, and Re-usable data³. Therefore, the security concept is
191 continuously extended to suggest best practices for data sharing, data publication, community features
192 on the platform, e.g. user forums, use of the summary feature, or use of social media within the platform.

193

194 **Description of data upload to FASTGenomics**

195 FASTGenomics allows registered users to upload own data sets in the data module for further analysis. In
196 the data upload window, the user is asked to provide the expression file in sparse format, the NCBI
197 taxonomy ID of the organism and a title for the later appearance in the dataset item list. Once the user

198 has provided the information, the uploaded data is transformed to our internal FASTGenomics Data
199 Package Format. A Python script checks the NCBI Entrez IDs and those gene IDs that cannot be mapped
200 uniquely. The script automatically documents the removed genes. Once finished, the data is provided to
201 the user in the data module. Thereby, the data set is by default only visible to the person who has
202 uploaded it.

203 In addition, we also provide the possibility to directly generate the FASTGenomics Data Package Format.
204 A documentation including an R-based tutorial for dataset preparation can be found at
205 https://github.com/FASTGenomics/FASTGenomics_Data_Package_Format. Here, the user can add
206 further information to the data set like metadata for the cells and genes and a dataset description
207 including a short abstract, and contact information. This FASTGenomics Data Package can then be
208 uploaded to the data module and allows fast access to the FASTGenomics functionalities.

209

210 **Description of Summary of any given analysis**

211 The analysis summary report gives a detailed overview of all steps executed and results produced in an
212 analysis workflow to facilitate understanding and to ensure reproducibility (**Supplementary Figure S5**).
213 Workflows describe the sequence of analysis steps performed to get from raw data to an analysis result.
214 Such workflows are not necessarily linear and may contain several branches (e.g. when one app depends
215 on input from several other apps executed before), i.e. a workflow is a directed acyclic graph. Each node
216 (i.e. a versioned calculation or visualization app) in this graph is described in the workflow definition,
217 including information on the analysis context within the workflow. Furthermore, each app provides text
218 passages containing information about applied methods and offers links to access generated interim
219 results during runtime. For report generation, the `summary_visualization` app recursively resolves
220 the app dependencies (required input data and necessary analysis steps to generate this) from leaf to the
221 root and dynamically assembles information gathered from the nodes into the final analysis report. Apps
222 connecting the analysis summary to laboratory information management systems (LIMS) will be
223 developed in the near future and included into the FASTGenomics analytical ecosystem.

224

225 **Concordance rate analysis of FASTGenomics Analyses with published results**

226 Analysis of single-cell RNA-seq data is a complex multistep procedure with many methods available for
227 individual tasks, however with no gold standard being defined. Published experiments thus typically
228 present analysis strategies that are highly specific for the respective underlying dataset. Accordingly,
229 comparison of analysis strategies can be a daunting task. Here, we applied the FASTGenomics subtype
230 discovery workflow for the analysis of a selection of published single-cell RNA-seq datasets generated with
231 different technologies and of various dataset sizes. One common task in single-cell RNA-seq analysis is the
232 definition of cell clusters to define sub-populations in complex mixtures of cells, with a definition of
233 characteristic gene expression signatures and their functional characterization being typical downstream
234 applications that crucially depend on the cluster assignment of cells. We therefore quantitatively

235 compared FASTGenomics single-cell cluster assignments based on a neural network-based dimensionality
236 reduction algorithm (see description below) to previously published clustering results for a selection of
237 single-cell RNA-seq datasets (**Supplementary Table 2**)⁴⁻⁹.

238

239 **Neural network-based dimensionality reduction and clustering**

240 The standard subtype discovery workflow in FASTGenomics consists of three calculation apps that reduce
241 the input dimensionality and group samples based on their similarity as seen in the gene expression
242 profile. The first calculation app normalizes the data using the term-frequency times inverse-document-
243 frequency (TF-IDF)¹⁰. This scheme replaces the gene expression in each sample with a number
244 proportional to the expression amplitude in this sample multiplied by the inverse number of samples in
245 which the gene is observed. This amplifies genes which are specific for a given subpopulation and
246 dampens the effect of genes which are present in most samples. Since the non-linear dimensionality
247 reduction needs a dense matrix with intermediate dimensionality, the sparse, normalized expression table
248 is compressed to 32 dimensions using truncated singular value decomposition implemented in a second
249 calculation app. In the third step, a calculation app uses a neural network to approximate a parametric t-
250 SNE embedding¹¹. This step projects the intermediate 32-dimensional data onto a two-dimensional space.
251 The neural network approximates the t-SNE optimization problem by learning a projection that minimizes
252 the t-SNE loss function and allows iterative training on batches. By default, the app uses batches consisting
253 of 512 samples and calculates the joint probabilities of samples in the higher dimensional space (by default
254 32 dimensions from the truncated singular value decomposition) and the target two-dimensional space.
255 Then, the network minimizes the Kullback-Leibler (KL) divergence between input probability and output
256 probability similar to the original t-SNE algorithm. Finally, clustering of the cells is performed using
257 the HDBSCAN algorithm.

258

259 **Methods applied to determine concordance**

260 ***Dataset pre-processing and clustering***

261 For cluster determination we downloaded count tables derived from previously published scRNA-seq
262 datasets (**Supplementary Table 2**) and used them as they were provided within the public repository.
263 After upload into FASTGenomics we used the pre-installed workflow 'subtype discovery' for the
264 identification of clusters within the dataset using the above described neural network-based
265 dimensionality reduction algorithm. Cluster assignments for individual cells were requested from the
266 corresponding authors of individual datasets and compared to those obtained by the FASTGenomics
267 workflow.

268 ***Concordance rate calculation***

269 From individual cells' published cluster assignments as well as clusterings produced with the
270 FASTGenomics subtype discovery workflow, a contingency matrix M of cell counts per cluster pairs was
271 generated. To provide a quantitative measure for the concordance of clustering results obtained from the

272 two methods, we calculated the concordance rate C for each cluster produced with a specific method as
273 follows:

$$274 \quad C(M_i) = \frac{\max(M_i, \cdot)}{\sum_j M_{i,j}} \times 100$$

275 To provide an overall summary statistic how well one clustering method captures results from the other
276 method, the median concordance rate was calculated. For each pair of clustering methods, two median
277 concordance rates (one for each method) can be calculated.

278 ***Adjusted mutual information***

279 To quantify the overall clustering concordance between published and the FASTGenomics analysis, we
280 use the adjusted mutual information (AMI), which ranges between 0 (two clusterings show only random
281 overlap) and 1 (overlap between two clusterings is not due to chance)⁷.

282

283 **Results of concordance rate analysis**

284 To provide an unbiased analysis of the publicly provided dataset we did not adjust the number of cells
285 when uploading the data to the FASTGenomics portal. As shown in **Supplementary Table 2**, the number
286 of cells reported in the respective publications and the number of cells publicly available was not always
287 identical. Due to the inclusion of different sets of cells into the analysis and the differences in analysis
288 settings, we are aware that the number of clusters between the FASTGenomics analysis and previously
289 published results can vary. Nevertheless, if publicly available datasets are to be used by a broader
290 community, we postulated that the use of the complete datasets provided will be the default usage of
291 such data resources.

292 We performed a quantitative comparison between previously reported cluster structures and clusters
293 determined in an unbiased fashion by the FASTGenomics ‘subtype discovery’ pipeline for a selection of
294 single-cell RNA-seq datasets (**Figure 2A**). We observed variation in the AMI values determined for the
295 selected datasets and hypothesize that apart from technical influences like sparsity of the expression
296 matrix and read coverage per cell, also biological aspects impact the clarity of cellular subtype discovery.
297 We found that AMI values were generally lower in immune cell datasets compared to those from other
298 cell types (cerebral and cancer cells as well as retinal tissue), presumably due to the lower RNA content of
299 immune cells and the lower number of genes expressed therein^{8,9}.

300 A more detailed analysis of the cluster structure was performed for 3,005 single-cell transcriptomes
301 derived from the murine primary somatosensory cortex (S1) and the hippocampal CA1 region⁴, which was
302 previously divided into 9 main clusters and 47 subclasses. When applying the FASTGenomics ‘subtype
303 discovery’ pipeline, we identified 16 clusters. Of the 3,005 cells analyzed in the published study⁴, 630 cells
304 (20.1%) could not be assigned to any of the clusters due to limited assignment confidence resulting from
305 almost equidistant positioning between cell clusters in the tSNE space. However, the median concordance
306 rates for the previously determined 9 main clusters and the 16 newly defined ones were as high as 96.5%

307 for FASTGenomics and 90% for BACKSPIN, arguing for a high degree of concordance for a large fraction of
308 clusters and cells. Likewise, the AMI for both clustering results was high ($AMI_{\text{Zeisel et al., Science (2015)}} = 0.75$,
309 **Figure 2A**). We further compared the 16 FASTGenomics clusters to the 7 cell classes corresponding to the
310 9 published main clusters as well as to the 47 cell subtypes (**Figure 2B**). Here, all oligodendrocytes classes
311 and all S1 and CA1 pyramidal neurons were each captured by one FASTGenomics cluster, while
312 interneurons were mainly represented in six distinct clusters by FASTGenomics. Thus, the standard
313 subtype discovery workflow revealed an intermediate resolution between the overall and the fine-grained
314 published analysis without generating contradictions to existing knowledge.

315

316 **Setup of ASAP, Granatum and SeqGeq for comparison with FASTGenomics**

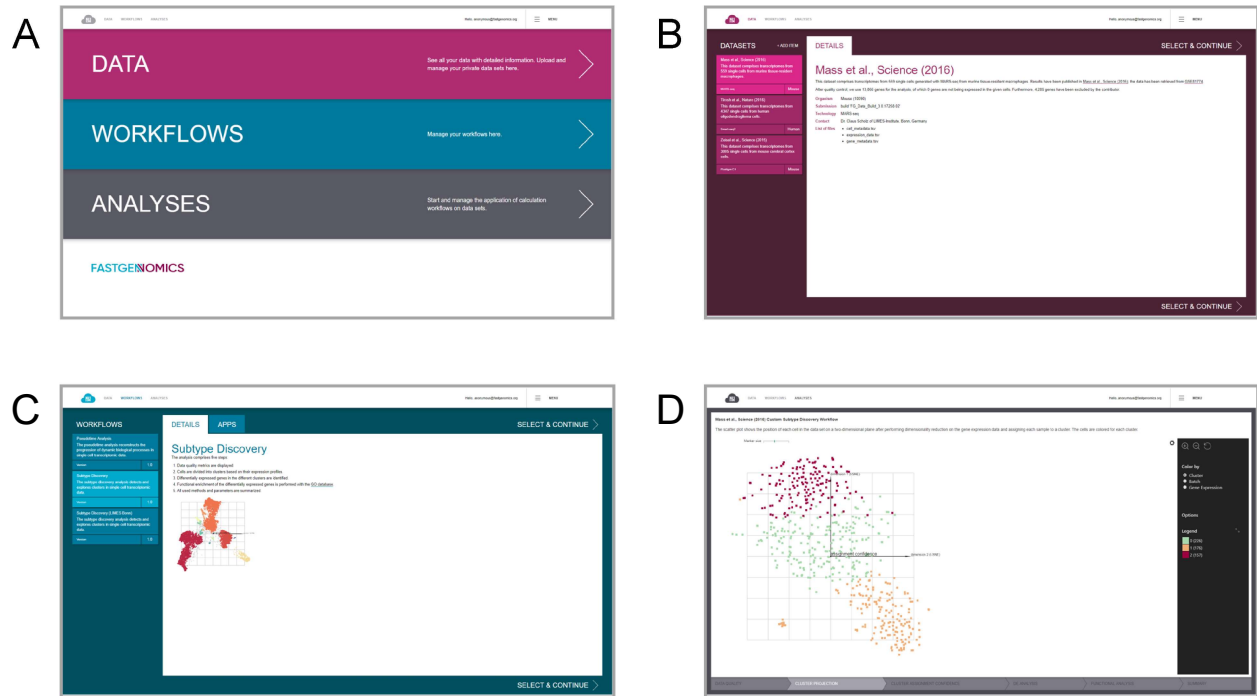
317 To evaluate the analytical capabilities of the FASTGenomics analysis pipeline, we defined a set of analytical
318 tasks and checked the performance on published single-cell analysis pipelines featuring a graphical user
319 interface, ASAP¹², Granatum¹³ and SeqGeq¹⁴. ASAP was evaluated using the publicly accessible online
320 instance at <https://asap.epfl.ch>. Granatum required the local installation of VirtualBox version 5.1.26
321 (Oracle) and the import of the Granatum appliance version 1.1_2 obtained from
322 <http://garmiregroup.org/granatum/app>. SeqGeq version 1.3 was obtained from
323 <https://www.flowjo.com/solutions/seqgeq> and installed locally in the default configuration. The
324 performance of Granatum and SeqGeq was evaluated on a 64-bit Windows 10 machine with Intel i7 6700K
325 CPU and 32 GB RAM.

326

327 **Resource Requirements of a FASTGenomics Analysis Workflow**

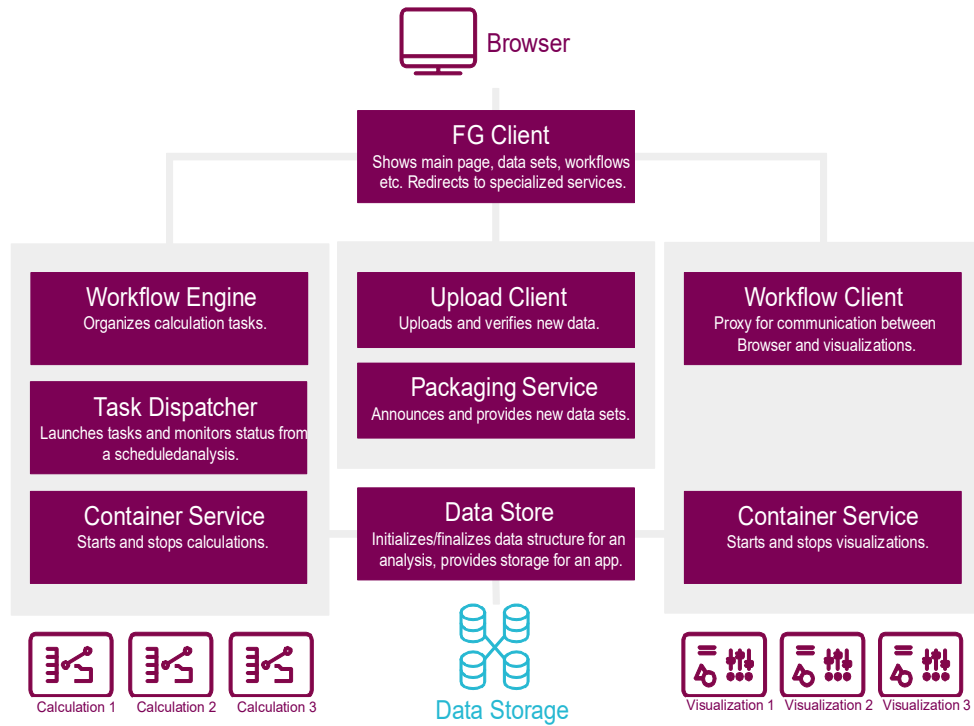
328 The memory requirements and the computing time to complete an analysis of a dataset of defined size
329 were chosen to describe the performance of the FASTGenomics pipeline. These parameters were
330 determined in a single-user setting for datasets consisting of 1,920 to 68,579 cells^{4,5,7-9,15}; analysis tasks
331 evaluated for the performance measurements were data normalization, dimensionality reduction and cell
332 clustering, because the effort needed for detection of differentially expressed genes and their functional
333 analysis depend on the number of clusters found in the single-cell dataset. Computing times for individual
334 analysis steps were extracted from the Docker log generated by the FASTGenomics Task Dispatcher and
335 summed up for all steps in the analysis workflow. Memory requirements were determined with a batch
336 script running in the background during the calculations that executes `docker stats` and `docker`
337 `ps` in intervals of two seconds logging memory consumption of individual containers. During the runtime
338 of each container, the maximum memory requirement was used for further evaluation. Resource
339 requirements were determined with the publicly accessible instance of FASTGenomics, which is currently
340 installed on a Standard D8s v3 system (8 vCPUs based on the 2.3 GHz Intel XEON[®] E5-2673 v4 processor
341 and 32 GB RAM).

342 **Supplementary Figures**
Figure S1



343
344 **Supplementary Figure 1: Graphical User Interface of the FASTGenomics analysis ecosystem. (A)** The
345 main menu. **(B)** Overview of accessible datasets. **(C)** Workflow selection page. **(D)** Analysis result
346 visualization.
347

Figure S2



348

349

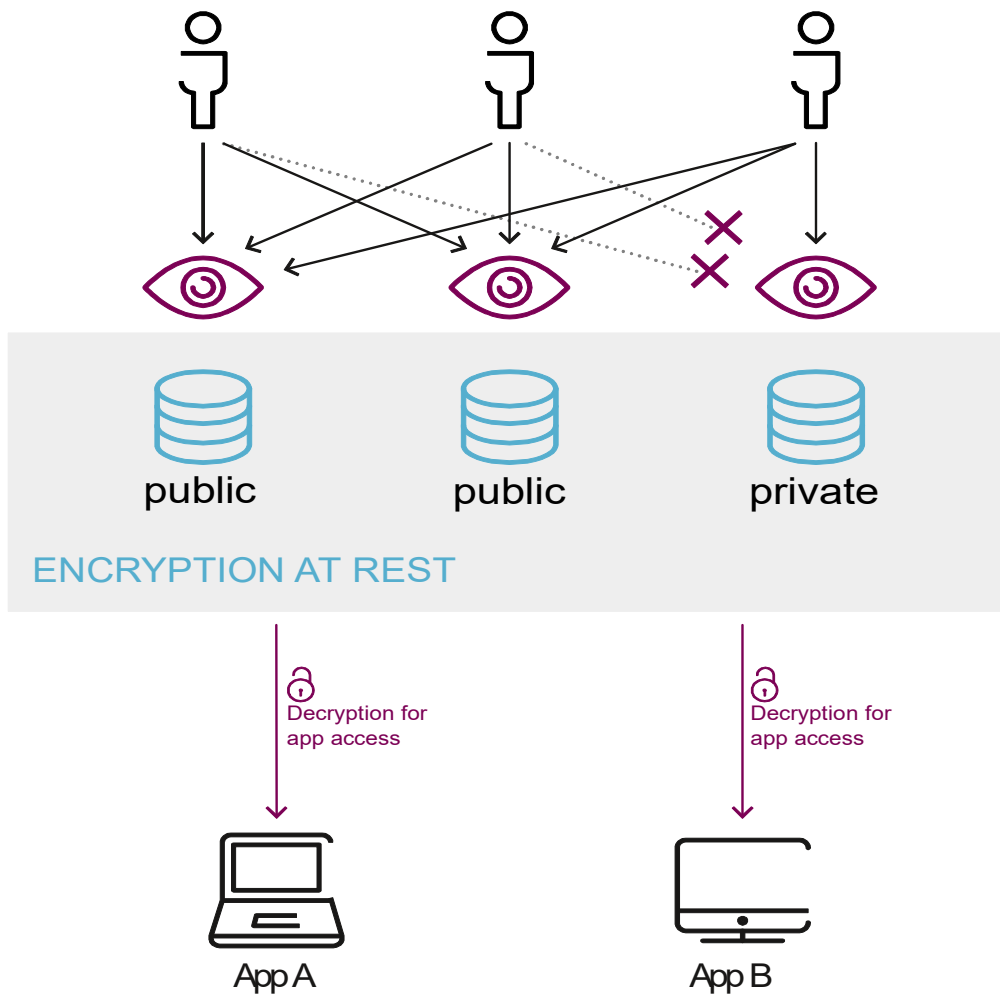
350

351

352

Supplementary Figure 2: Technical Representation of the FASTGenomics Architecture. A web browser allows the interaction with the FASTGenomics Client, which internally manages the upload dock (middle branch), the calculation engine (left branch) and the screenflow engine (right branch).

Figure S3



353
354
355
356
357

Supplementary Figure 3: Data Access. All data is stored on encrypted storage devices. Users can access public datasets and those uploaded by the users, but not those from other users. Data is encrypted for calculation and visualization apps.

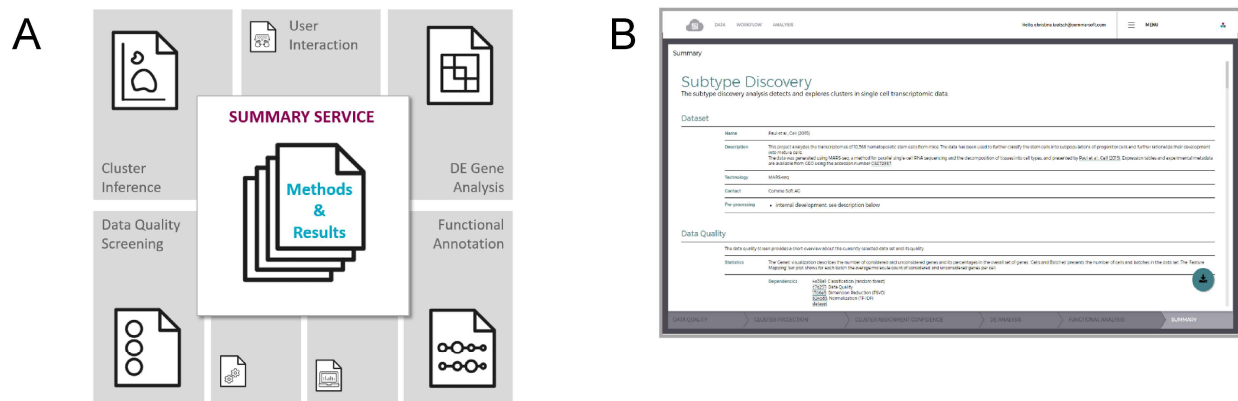
Figure S4



358
 359 **Supplementary Figure 4: Data Quality. (A)** Data upload in FASTGenomics. **(B)** Data quality check concept.
 360 Based on the overall number of genes detected in the single-cell transcriptomics dataset, a lower
 361 threshold for the number of genes per cell is dynamically generated; cells expressing fewer genes are
 362 excluded from analysis. Furthermore, genes expressed in less than a predefined proportion of cells are
 363 removed from the dataset. **(C-E) Quality check screens in FASTGenomics**, screenshots illustrating **(C)**
 364 average molecule counts, **(D)** gene types, **(E)** quantification of batch effects.

365

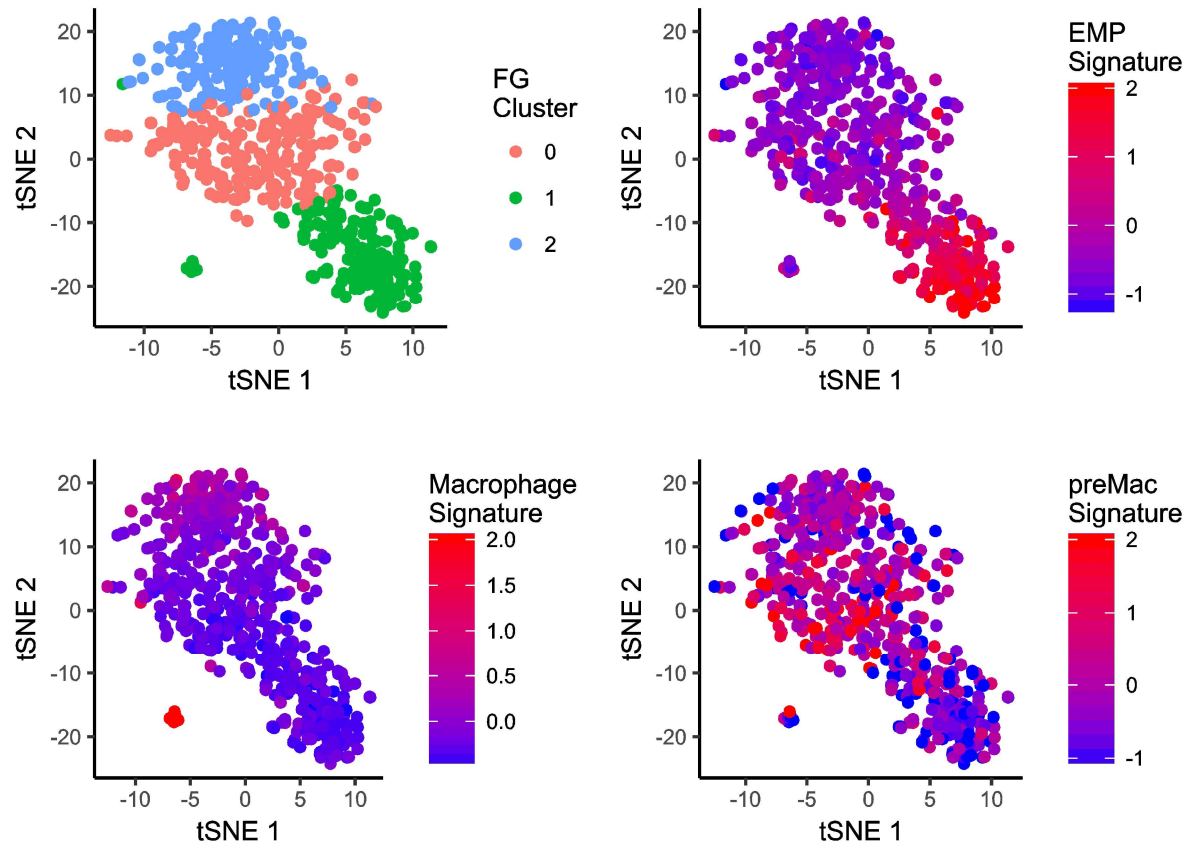
Figure S5



366
367 **Supplementary Figure 5: Analysis Summary.** (A) Schematic overview of the summary report-generating
368 app. Information about analyzed data, applied workflows with invoked apps including version information
369 and parametrization are recursively resolved and compiled during run-time to generate the summary
370 report. (B) Screenshot of an example summary report generated during the subtype discovery workflow.

371

Figure S6



372

373 **Supplementary Figure 6: Re-Analysis of Mass et al. (A)** Cluster assignments of individual cells. **(B)**

374 Individuals cells colored according to EMP signature gene expression. **(C)** Individual cells colored

375 according to macrophage signature gene expression. **(D)** Individual cells colored according to preMac signature gene

376 expression.

377

| Name | Functionality | Algorithm/Method | Status |
|-------------------------------|---|----------------------------|-------------------|
| calc_batch_effect_classifier | Batch effect quantification | Random Forest | Active |
| calc_clustering_hdbscan | Cell clustering | HDBSCAN | Active |
| calc_clustering_louvain | Cell clustering | Louvain | Under Development |
| calc_count_normalize | Normalization | various | Under Development |
| calc_de_genes_diffrank | Detection of differentially expressed genes | Diffrank (Scanpy) | Active |
| calc_de_genes_glm | Detection of differentially expressed genes | Generalized Linear Model | Active |
| calc_de_genes_nonparametric | Detection of differentially expressed genes | Mann-Whitney U test | Active |
| calc_diffusion_pseudotime | Pseudo-temporal ordering of cells | Scanpy | Active |
| calc_dimreduction_autoencoder | Dimensionality reduction | Neural network | Active |
| calc_dimreduction_tsne | Dimensionality reduction | tSNE | Active |
| calc_filter_quality | Quality control | Detection rate filtering | Under Development |
| calc_functional_analysis | Enrichment analysis | Fisher's Exact test | Active |
| calc_list_filtering | Blacklist/whitelist filtering | ID list filtering | Under Development |
| calc_logreg_confusion | | Logistic regression | Active |
| calc_normalize_tfidf | Normalization | TF/IDF | Active |
| calc_tsvd | Dimensionality reduction | Single value decomposition | Active |

378 **Supplementary Table 1A:** Calculation Apps in FASTGenomics. The table lists the calculation apps
379 available in FASTGenomics and currently under development. For each app, the name, functionality and
380 used algorithm or method is specified.

| Name | Usage | Status |
|---------------------|-----------------------------------|---|
| viz_barchart | Data Quality | Active |
| viz_batch_effect | Data Quality | Active |
| viz_dataquality | Data Quality | Active |
| viz_confusionmatrix | Cluster Inference | Active |
| viz_heatmap | DE Genes, Functional Analysis | Active |
| viz_scatterplot | Clustering, Diffusion Pseudo-time | Active |
| viz_table | DE Genes | Active |
| viz_linechart | Gene dynamics | Under development/ soon in FAST Genomics |

381 **Supplementary Table 1B:** Visualization Apps in FASTGenomics. The table lists the visualization apps
382 available in FASTGenomics and currently under development.

| ID | GEO Accession Code | Organism | Number of Genes | Number of Cells |
|--|---------------------------|-----------------|------------------------|------------------------|
| Macosko et al., Cell (2015) | GSE63473 | Mouse | 21,605 | 49,300 |
| Mass et al., Science (2016) | GSE81774 | Mouse | 8,553 | 408 |
| Moignard et al., Nature Biotechnology (2015) | --- | Mouse | 46 | 3,934 |
| Nestorowa et al., Blood (2016) | GSE81682 | Mouse | 23,357 | 1,920 |
| Paul et al., Cell (2015) | GSE72857 | Mouse | 19,362 | 10,368 |
| Tirosh et al., Nature (2016) | GSE70630 | Human | 22,338 | 4,347 |
| Tirosh et al., Science (2016) | GSE72056 | Human | 22,333 | 4,645 |
| Zeisel et al., Science (2015) | GSE60361 | Mouse | 18,920 | 3,005 |
| Zheng et al., Nature Communications (2017) | --- | Human | 21,253 | 68,579 |
| Ziegenhain et al., Molecular Cell (2017) | GSE75790 | Mouse | 22,701 | 482 |

383 **Supplementary Table 2:** Description of datasets available in the FASTGenomics pipeline.

384 **Supplementary References**

- 385 1. Internetauftritt der Bundesbeauftragten für den Datenschutz und die Informationsfreiheit -
386 Homepage - Bundesdatenschutzgesetz (BDSG). Available at:
387 <https://www.bfdi.bund.de/SharedDocs/Publikationen/GesetzeVerordnungen/BDSG.html>. (Accessed:
388 8th November 2017)
- 389 2. European Parliament and the Council of the European Union. Regulation (EU) 2016/679 of the
390 European Union Parliament and of the Council of 27 April 2016. (2016).
- 391 3. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and stewardship.
392 *Sci. Data* **3**, 160018 (2016).
- 393 4. Zeisel, A. *et al.* Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-
394 cell RNA-seq. *Science* **347**, 1138–1142 (2015).
- 395 5. Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual Cells Using
396 Nanoliter Droplets. *Cell* **161**, 1202–1214 (2015).
- 397 6. Mass, E. *et al.* Specification of tissue-resident macrophages during organogenesis. *Science* **353**,
398 (2016).
- 399 7. Nestorowa, S. *et al.* A single-cell resolution map of mouse hematopoietic stem and progenitor cell
400 differentiation. *Blood* **128**, e20-31 (2016).
- 401 8. Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq.
402 *Science* **352**, 189–196 (2016).
- 403 9. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* **8**,
404 14049 (2017).
- 405 10. Jones, K. S. A statistical interpretation of term specificity and its application in retrieval. *J. Doc.*
406 **28**, 11–21 (1972).

- 407 11. Maaten, L. Learning a Parametric Embedding by Preserving Local Structure. in *PMLR* 384–391
408 (2009).
- 409 12. Gardeux, V., David, F. P. A., Shajkofci, A., Schwalie, P. C. & Deplancke, B. ASAP: a Web-based
410 platform for the analysis and interactive visualization of single-cell RNA-seq data. *Bioinforma. Oxf. Engl.* (2017). doi:10.1093/bioinformatics/btx337
- 412 13. Zhu, X. *et al.* Granatum: a graphical single-cell RNA-Seq analysis pipeline for genomics scientists.
413 *Genome Med.* **9**, 108 (2017).
- 414 14. SeqGeq® | FlowJo, LLC. Available at: <https://www.flowjo.com/solutions/seqgeq>. (Accessed: 2nd
415 February 2018)
- 416 15. Paul, F. *et al.* Transcriptional Heterogeneity and Lineage Commitment in Myeloid Progenitors.
417 *Cell* **163**, 1663–1677 (2015).
- 418