

1 **The Male Fertility Gene Atlas – A web tool for collecting and integrating data about**  
2 **epi-/genetic causes of male infertility**

3

4 Running title: Male Fertility Gene Atlas (MFGA)

5

6 H. Krenz<sup>1</sup>, J. Gromoll<sup>2</sup>, T. Darde<sup>3</sup>, F. Chalmel<sup>3</sup>, M. Dugas<sup>1</sup>, F. Tüttelmann<sup>4\*</sup>

7

8 <sup>1</sup>Institute of Medical Informatics, University of Münster, 48149 Münster, Germany

9 <sup>2</sup>Centre of Reproductive Medicine and Andrology, University Hospital Münster, 48149 Münster,  
10 Germany

11 <sup>3</sup>Univ Rennes, Inserm, EHESP, Irset (Institut de recherche en santé, environnement et travail)  
12 - UMR\_S 1085, F-35000 Rennes, France

13 <sup>4</sup>Institute of Human Genetics, University of Münster, 48149 Münster, Germany

14

15 \*Correspondence address: Institute of Human Genetics, University of Münster, Vesaliusweg  
16 12-14, 48149 Münster, Germany, [frank.tuettelmann@ukmuenster.de](mailto:frank.tuettelmann@ukmuenster.de)

17

18

19 Word count: abstract 276, main text: 3914

20

21 **Abstract**

22 Interconnecting results of previous OMICs studies is of major importance for identifying novel  
23 underlying causes of male infertility. To date, information can be accessed mainly through  
24 literature search engines and raw data repositories. However, both have limited capacity in  
25 identifying relevant publications based on aggregated research results e.g. genes mentioned  
26 in images and supplements. To address this gap, we present the Male Fertility Gene Atlas  
27 (MFGA), a web tool that enables standardised representation and search of aggregated result  
28 data of scientific publications. An advanced search function is provided for querying research  
29 results based on study conditions/phenotypes, meta information and genes returning the exact  
30 tables and figures from the publications fitting the search request as well as a list of most  
31 frequently investigated genes. As basic prerequisite, a flexible data model that can  
32 accommodate and structure a very broad range of meta information, data tables and images  
33 was designed and implemented for the system. The first version of the system is published at  
34 the URL <https://mfga.uni-muenster.de> and contains a set of 46 representative publications.  
35 Currently, study data for 28 different tissue types, 32 different cell types and 20 conditions is  
36 available. Also, ~5,000 distinct genes have been found to be mentioned in at least ten of the  
37 publications. As a result, the MFGA is a valuable addition to available tools for research on the  
38 epi-/genetics of male infertility. The MFGA enables a more targeted search and interpretation  
39 of OMICs data on male infertility and germ cells in the context of relevant publications.  
40 Moreover, its capacity for aggregation allows for meta-analyses and data mining with the  
41 potential to reveal novel insights into male infertility based on available data.

42

43

44 **Keywords:** male infertility, genetics, epigenetics, omics, database

45

46 **Introduction**

47 Male infertility is a prevalent and highly heterogeneous disease for which the underlying causes  
48 can currently be identified for only about 30% of male partners in infertile couples (Tüttelmann  
49 *et al.*, 2018). In the last years, a growing number of studies have been published on the  
50 genetics and epigenetics of male infertility using a broad range of genomic technologies. For  
51 an overview see Oud *et al.* (2019). Often, researchers give extensive insight into their findings  
52 by providing supplementary data or access to their raw data. Theoretically, this enables  
53 clinicians and other researchers to validate those findings or interpret their own results in a  
54 broader context; in practice, however, the findability and reusability of this vast set of  
55 information is limited. To date, there is no public resource for the field of male infertility that  
56 enables clinicians and researchers alike to easily access a comprehensive overview of recent  
57 findings, although tools for other diseases have long been established, such as, for example,  
58 the Gene ATLAS (Canela-Xandri *et al.*, 2018) that was released in 2017. It provides information  
59 on genetic associations of 778 traits identified in genome-wide association studies (Canela-  
60 Xandri *et al.*, 2018), male infertility is, unfortunately, not one of them.

61 Instead of being accessible through a specified tool, research on male infertility relies mainly  
62 on general search engines for scientific publications like PubMed or Google Scholar. These  
63 engines can be employed to search for information based on keywords of interest, e.g., gene  
64 names in combination with specific conditions. However, if the required information can only  
65 be found in a supplementary table or figure, these search engines are unable to mark the  
66 corresponding publication as relevant. Other sources of information are raw data repositories  
67 like Gene Expression Omnibus (GEO) which provides access to the raw data files of  
68 microarray and genomic data from many publications (Barrett *et al.*, 2013), or Sequence Read  
69 Archive (SRA). This is a marked achievement for the scientific community, but such  
70 repositories do not allow users to find a data file based on a gene or variant of interest. For the  
71 most part, publications of interest have to be identified in advance, and comprehensible insight  
72 into the data can only be achieved by reanalysing it, using complex bioinformatics pipelines.

73 To address this need specifically for the reproductive sciences, Chalmel and colleagues  
74 established the ReproGenomics Viewer (RGV) in 2015 (Darde *et al.*, 2019; Darde *et al.*, 2015).  
75 It provides a valuable resource of manually-curated transcriptome and epigenome data sets  
76 processed with a standardised pipeline. RGV enables the visualisation of multiple data sets in  
77 an interactive online genomics viewer and, thus, allows for comparisons across publications,  
78 technologies and species (Darde *et al.*, 2019). However, the RGV is not designed to show  
79 downstream analysis results such as differentially expressed genes or genomic variation. Also,  
80 it is not possible to identify relevant publications based on genes of interest. Other tools in this  
81 field are GermOnline (Lardenois *et al.*, 2010) and SpermatogenesisOnline (Zhang *et al.*, 2013).  
82 However, both have not been maintained in recent years.

83 In order to offer a public platform that provides access to a comprehensive overview of  
84 research results in the field of epi-/genetics of male infertility and germ cells and to bridge the  
85 gap between textual information in publications on the one hand and complex data sets on the  
86 other hand, we designed, developed, and now introduce the publicly available Male Fertility  
87 Gene Atlas (MFGA, <https://mfga.uni-muenster.de>). Its objective is to provide fast, simple and  
88 straightforward access to aggregated analysis results of relevant publications, namely by  
89 answering questions like “What is known about the gene *STAG3* in the context of male  
90 infertility?” or “Which genes have been identified to be associated with azoospermia?”. To this  
91 end, we created an advanced search interface as well as comprehensive overviews and  
92 visualisations of the publications and search results. A basic prerequisite was that we had to  
93 design and implement a data model that can accommodate and structure a very broad range  
94 of meta-information, data tables and images.

95

## 96 **Materials and methods**

### 97 **Requirements engineering**

98 The MFGA is designed to support researchers and clinicians in the field of male infertility in  
99 finding and evaluating the analysis results of relevant publications. To ensure that the MFGA

100 offers a relevant scope and focusses on the central needs and interests of the targeted user  
101 group, an extensive requirements engineering process had to be employed. This process  
102 aimed at defining the most relevant parts of the system design; in case of the MFGA, this  
103 means answering the following five questions:

- 104 1. What kind of data should be available?
- 105 2. How should the database be maintained and kept up to date?
- 106 3. Who should be authorised to access the data on the website?
- 107 4. What interface should be provided to access the data?
- 108 5. How should the data be visualised?

109 Since user-centric design is proven to be a critical factor of success for software projects  
110 (Maguire and Bevan, 2002), the requirements were acquired in close cooperation with a large  
111 group of prospective users of the MFGA. These researchers and clinicians are part of the  
112 Münster-based clinical research unit (CRU326) (<http://www.male-germ-cells.de>) and provided  
113 input on a broad range of aspects of male infertility, germ cell development as well as sperm  
114 morphology, motility, and function. The development of the MFGA is based on the software  
115 development and requirements engineering paradigm Rapid Prototyping, which focusses on  
116 developing early testing versions of the product to enable informed user feedback and iterative  
117 improvements (Budde *et al.*, 1992; Gordon and Bieman, 1995).

118 The requirements analysis was based on a set of 39 representative and relevant  
119 publications/data sets that the MFGA should be able to host, provided by members of the  
120 CRU326. The publications were highly heterogeneous and covered a broad range of  
121 methodologies, from GWA studies to targeted genotyping, microarray, bulk or single-cell RNA  
122 sequencing as well as methylation. The complete list is available in Suppl. Tab. S1 and has  
123 already been included into the MFGA database.

124

125 **Current state of requirements**

126 The requirements regarding the kind of data that should be available in the MFGA can be  
127 summarised in two points: First, the MFGA should host data representing the analysis results  
128 of a comprehensive set of publications on male infertility. Generally speaking, the MFGA needs  
129 to be able to deal with publications of arbitrary form and content. To account for the extremely  
130 high degree of heterogeneity between publications, a very flexible, modular data model is  
131 required. However, to enable the MFGA's central service of offering selective data retrieval  
132 operations, it is crucial to identify, standardise and tag recurrent relevant information items and  
133 to curate them manually for each publication. Whenever applicable, standardisation should be  
134 based on ontologies as provided by OBO Foundry (Smith *et al.*, 2007). Second, the MFGA  
135 should provide complementary data from external databases to increase the comprehensibility  
136 of information, e.g., to explain the background of gene names and provide links to more details.  
137 However, the MFGA is not supposed to contain any raw data; instead the system should link  
138 to the corresponding repositories if provided by a publication. Also, the MFGA is not planned  
139 for processing raw data by itself; the system is only intended to show aggregated result data  
140 provided by publications.

141 The database of the MFGA should be fully public and accessible to all researchers and  
142 clinicians in the field of male infertility. Maintenance and curation of selected publications  
143 should be provided by the development team. In order to keep the database up to date in a  
144 structured way, a recurring process for preparing new relevant publications will be  
145 implemented (Fig. 1). It comprises three main steps: selection of relevant publications, manual  
146 curation according to the requirements of the MFGA database and data upload. Prospectively,  
147 it should be triggered once per month and combine the proposed publications of four different  
148 sources into one prioritised list. The four information sources are (1) the list of not yet  
149 processed publications from the previous month, (2) newly published manuscripts from the  
150 CRU326, (3) other relevant new publications identified by the MFGA team, and, as an  
151 additional unbiased source, (4) recent publications queried from PubMed (query: "(male OR  
152 men) AND (fertility OR infertility) AND (gene OR genetic) AND ("YYYY/MM/01"[Date -

153 Publication] : "YEAR/MM/31"[Date - Publication]"). Publications will be curated and uploaded  
154 by priority. Eventually, users are supposed to be enabled to act as registered contributors and  
155 support selection, curation and upload of content. Thus, a user management and  
156 authentication system is required.

157 Regarding the interface for accessing the information, users should be able to view a complete  
158 list of publications and their key features to get a quick overview on the full content of the data  
159 base. Additionally, an advanced search function should provide identifying and showing  
160 subsets of relevant publications based on the occurrence of gene names or IDs in tables and  
161 figures as well as meta-information like data type, OMICS, processed tissues/cells and  
162 conditions/species. Detailed information and results of each individual publication should be  
163 displayed by comprehensible, standardised overviews, including images, tables and on-  
164 demand plots. As a special feature, the MFGA should display data from different publications  
165 in an aggregated way to allow for meta-analyses and integrated analyses. All pages should be  
166 enriched with appropriate additional information from external databases and with cross  
167 references, allowing for quick navigation through the atlas.

168

## 169 **Results**

### 170 **IT architecture**

171 The MFGA has been implemented as a modern *Java* web application that is securely hosted  
172 on a server at the University of Münster, Germany, and can be accessed freely via the URL  
173 <https://mfqa.uni-muenster.de>. Its architecture (Suppl. Fig. S1) is based on the specific  
174 requirements, identified during requirements engineering and explained in detail in the  
175 supplement. For a good user experience, the graphical user interface utilises the *Bootstrap*  
176 library which enables a smooth adjustment of the application to different browsers and devices  
177 (Otto *et al.*, 2020) and familiar design elements. On the server side, Spring Security framework  
178 is employed to provide secure authentication processes and manage data access privileges  
179 (Pivotal Software Inc., 2019).

180 **Data model**

181 The specific relational data model of the MFGA as well as the preparation process for  
182 publication data are a direct result of the requirements analysis and were developed based on  
183 the initial set of relevant publications. They are the main prerequisite for enabling complex  
184 search queries, aggregation functions and meta analyses on the MFGA's database. For  
185 representing the highly heterogeneous publications, a flexible and thus modular data model  
186 was designed (Suppl. Fig. S2). It is based on information items that were identified to be  
187 preserved over large subsets of publications on male infertility and standardised to seven  
188 classes of information items: *publication meta information*, *data set meta information*,  
189 *processed tissue type*, *processed cell type*, *cohort's condition*, *image* and *table* (Tab. 1). For  
190 each publication, they can freely be combined in any quantity including zero. In order to prevent  
191 inconsistencies in the database caused by misspelling or synonymous terms, structured data  
192 entry with drop-down lists and ontologies is implemented. Currently, the MFGA employs  
193 BRENDA Tissue Ontology (BTO) (Gremse *et al.*, 2011) for tissue type definition, Cell Ontology  
194 (CL) (Diehl *et al.*, 2016) for cell types and Human Phenotype Ontology (HPO) (Köhler *et al.*,  
195 2019) for conditions. Further technical specification, especially representation and annotation  
196 of data tables and images, is explained in the extended methods (see supplement).

197

198 **Available content**

199 Currently there are 46 publications available on the MFGA, and they comprise information on  
200 101 data sets (Fig. 2), 28 different tissue types, 32 different cell types and 20 conditions. The  
201 majority of data sets contains data on the transcriptome (22 single-cell RNA sequencing, 12  
202 bulk RNA sequencing & 3 microarray). The second largest group of data sets contains  
203 information on the genome / exome (13 targeted genotyping, 9 genome-wide association  
204 studies, 4 whole-genome sequencing & 4 sanger sequencing). Additionally, 6 data sets on  
205 whole-genome bisulfite sequencing and 2 on targeted deep bisulfite sequencing are available.



206 The most frequent conditions in the data sets are variants of abnormal spermatogenesis, e.g.  
207 azoospermia and oligozoospermia. Testicular tissue has been processed by 32 proteome,  
208 methylome and transcriptome data sets. The predominantly represented cell types are  
209 embryonic stem cells, primordial germ cells, spermatogonia, spermatocytes, and spermatids.  
210 Also, ~5,000 distinct genes have been found to be mentioned in at least ten of the publications,  
211 with the top genes being *TEX11*, *TEX14*, *TEX15*, *SYCP3*, *SMC1B*, *REC8* and *HORMAD1*.  
212 The full list of 46 publications can be accessed via the URL [https://mfqa.uni-](https://mfqa.uni-muenster.de/publication.html)  
213 [muenster.de/publication.html](https://mfqa.uni-muenster.de/publication.html).

214

## 215 **Functionality**

216 The MFGA provides a web interface for fast, simple, and straightforward access to the results  
217 of publications on the epi-/genetics of male infertility and germ cells. For this purpose, an  
218 advanced search form is provided on the *Search* tab of the atlas (Suppl. Fig. S3), enabling the  
219 specification of search requests for relevant data sets. This function supports search terms  
220 from seven categories: OMICS, data type (e.g., single-cell RNA sequencing or targeted  
221 sequencing), condition, species, cell type, tissue type and gene/ID. For each of the first six  
222 categories, one can choose from a list of search terms equivalent to the information items  
223 stored in the database. Condition, cell type and tissue type are based on appropriate  
224 ontologies: HPO (Köhler *et al.*, 2019), CL (Diehl *et al.*, 2016) and BTO (Gremse *et al.*, 2011).  
225 The number of available data sets is shown in brackets after each search term. Gene names  
226 and IDs can be specified in a text field as a comma-separated list. Also, there are three search  
227 modalities: (1) users can perform a broad search to identify all data sets that are annotated  
228 with at least one of the specified search terms, (2) users can perform a more targeted search  
229 for all data sets that are annotated with at least one of the specified search terms per category  
230 and (3) users can perform a very specific search for all data sets that are annotated with each  
231 of the search terms (default option). As an example, Fig. 3 shows an extract of the output of a  
232 simple search request for data sets in the MFGA database containing information on the gene  
233 *STAG3* (Suppl. Fig S2 for the full screenshot). This refers back to the introductory question:

234 “What is known about the gene *STAG3* in the context of male infertility?”. For the following  
235 examples, the gene *STAG3*, which is located on chromosome 7, encodes a protein involved  
236 in meiosis and can be related to male infertility (van der Bijl *et al.*, 2019), is employed to explain  
237 the general functionalities of the MFGA. A detailed *Walk Through* is provided on  
238 <https://mfga.uni-muenster.de/walkThrough.html>.

239 Executing a search request on the MFGA results in a list of data sets matching the  
240 requirements, represented via charts and result tables (Fig. 3 and Suppl. Fig. S3). The left  
241 chart shows the number of returned data sets grouped by publication. Multiple data sets  
242 matching the search request in one publication can indicate that the authors provided further  
243 proof for their findings, such as e.g., van der Bijl *et al.* (2019), who screened two cohorts of  
244 infertile men. A textual summary of the data sets returned by the search request and some  
245 important meta-information is given by the *Data sets* table. All tables in the MFGA format are  
246 fully interactive such that the columns can be sorted, filtered and plotted online. Whenever the  
247 search contains genes or IDs, the right chart represents their frequencies in the different  
248 publications. This provides an initial estimation of relevance. Additionally, a second table lists  
249 their specific occurrences in data tables and figures. In the case when a search request is  
250 targeted at revealing relevant genes, such as, e.g., in the introductory question: “Which genes  
251 have been identified to be associated with azoospermia?” (Suppl. Fig. S4), the right chart and  
252 second table present the genes that occur most frequently in data tables of the returned data  
253 sets.

254 The search results returned by the MFGA are designed for quick navigation and information  
255 access. Therefore, many cross references are provided: Table entries directly link to overview  
256 pages of the corresponding publication and its data sets. Data tables and images are linked in  
257 the MFGA as well. However, this functionality is restricted to publications that are published  
258 under an open access license. As an example, Fig. 3 shows that the publication van der Bijl *et al.*  
259 *et al.* (2019) mentions *STAG3*. Clicking on the data set title leads to the overview page of that  
260 publication (Suppl. Fig. S5). Also, data tables containing the searched gene can directly be  
261 accessed in MFGA by a single click, e.g., supplementary table 3 of van der Bijl *et al.* (2019)

262 which is shown partly in Fig. 4 (full version in Suppl. Fig. S6). For a deeper analysis of the  
263 underlying read data of individual data sets, the MFGA provides a link to the RGV, whenever  
264 a data set is available in both tools, such as, e.g., Hammoud *et al.* (2015) in Fig. 3.

265 Data from external public databases has been integrated in order to enrich information on the  
266 publications shown in the MFGA. Throughout the application, gene names can be selected in  
267 table cells and plots to open up an overlay with further explanations (for an example, see  
268 Fig. 5). The overlay bundles textual information from RefSeq (O'Leary *et al.*, 2016) and HGNC  
269 (Yates *et al.*, 2017; HGNC Database, 2018). Additionally, a link to the corresponding  
270 GeneCards page (Stelzer *et al.*, 2016) is provided, as well as a shortcut to the MFGA search  
271 for that gene. Data from the GTEx project (Lonsdale *et al.*, 2013; The Broad Institute of MIT  
272 and Harvard, 2019) is employed to show the gene's expression in different tissues scaled by  
273 the largest median. The expression of the gene in testis tissue is highlighted in red.

274

## 275 **Discussion**

276 The MFGA provides a public platform to support researchers and clinicians in obtaining an  
277 overview of the recent findings in the field of male infertility and germ cells. The web-based  
278 tool includes recent libraries and methods for an improved user experience and straightforward  
279 usability. In order to enable an advanced search for relevant publications, a relational data  
280 model has been developed and implemented for knowledge representation. It records the  
281 recurrent information items in a structured and consistent way and can accommodate arbitrary  
282 publications from the field of male infertility, regardless of the kind of data analysis and  
283 technology. The search form enables a broad range of simple to very complex search queries  
284 such as "What is known about the gene *STAG3* in the context of male infertility?", "Which  
285 genes have been identified to be associated with azoospermia?" or "Which genes have been  
286 found to be expressed in Sertoli cells of human testis tissue using single cell RNA  
287 sequencing?" (Suppl. Fig. S7).

## Male Fertility Gene Atlas (MFGA)

---

288 The main purpose of the MFGA is to enable researchers and clinicians to consider their own  
289 research results in the context of other relevant literature. To this end, the system enables a  
290 highly selective search for studies with comparable conditions and parameters and offers the  
291 means to quickly review their main analysis results. Additionally, searches for genes can be  
292 performed in order to identify data sets containing the corresponding genes in their analysis  
293 results. This functionality supports the validation of candidate genes. Further, supplementary  
294 information from various sources is embedded into the MFGA, e.g., RefSeq (O'Leary *et al.*,  
295 2016), HGNC (Yates *et al.*, 2017; HGNC Database, 2018) and GTEx project (Lonsdale *et al.*,  
296 2013; The Broad Institute of MIT and Harvard, 2019).

297 Compared to general search engines like Google Scholar and PubMed, with millions of  
298 records, the MFGA will always return smaller numbers of search results. However, in the  
299 domain of male infertility, the relevance of the individual results in relation to the corresponding  
300 search terms used in the MFGA is expected to be greater than when using those same search  
301 terms in a large search engine. Search engines for literature are usually restricted to mining  
302 textual information of publications and cannot consider information that is presented in tables,  
303 images, or supplementary material. The MFGA, however, is specifically designed for that task,  
304 which is enabled by its key features: A standardised data representation and disease-specific  
305 manual curation.

306 Another problem occurs when relevant data sets are searched in raw data repositories like  
307 GEO or processed data resources like RGV (Darde *et al.*, 2015; Darde *et al.*, 2019). While the  
308 information that these tools are able to provide is much more detailed than overviews in the  
309 MFGA, they cannot be used to identify a data set based on, for example, a list of genes that  
310 are supposed to be differentially expressed, since raw data does usually not include that kind  
311 of annotation. The MFGA facilitates searching through aggregated result data and, thus,  
312 identifying relevant data sets. Once identified, such data sets might then be further investigated  
313 in the RGV (Darde *et al.*, 2015; Darde *et al.*, 2019) or by processing GEO data with  
314 bioinformatics tools. Since RGV (Darde *et al.*, 2015; Darde *et al.*, 2019) and MFGA

315 complement each other in their functionality, the MFGA links to RGV whenever a data set is  
316 present in both tools; further, an even closer integration is planned.

317 Regarding tools that might appear similar to the MFGA, Zhang *et al.* (2013) provide a tool  
318 termed SpermatogenesisOnline for searching individual genes in the context of  
319 spermatogenesis based on a set of publications from 2012 and earlier. However, it remains  
320 unclear which publications were included and whether or not they have been updated since  
321 then. Another tool from Lardenois *et al.* (2010), GermOnline, provides functional information  
322 about individual genes and access to transcriptomics data from microarrays on germline  
323 development from 2008. Both tools have, to our knowledge, not been updated and are, thus,  
324 not very useful anymore.

325 There are some limitations to the approach of the MFGA. Since analysis results are not  
326 reproduced using in-house pipelines, the MFGA is restricted to the analysis results authors are  
327 presenting in their publications. In the case where a publication provides, e.g., only significant  
328 SNPs from a GWA study, the MFGA, too, reports only these. The only meta-analyses enabled  
329 are those that use the aggregation of information in the MFGA, e.g., the top score of gene  
330 appearances, to approximate gene importance; the information cannot indicate correlation or  
331 even causality. Thus, the MFGA focusses on proposing plausible research hypotheses.  
332 Finally, full functionality of the MFGA can only be provided for publications under an open  
333 access license. Publications that allow no or restricted reuse are only partly integrated.

334 Prospectively, the MFGA will be updated and enlarged based on the proposed content  
335 management process (Fig. 1). In this context, a web form will be implemented to enable users  
336 to propose publications that are, from their point of view, still missing. Additionally, a closer  
337 integration of RGV and MFGA is planned, e.g. by coupling the publication submission forms,  
338 as well as further tools for meta-analyses.

339 In conclusion, the MFGA is a valuable addition to available tools for research on the epi-  
340 /genetics of male infertility. It helps fill the gap between pure literature searches as provided  
341 by PubMed and raw data repositories like GEO or SRA. The MFGA enables a more targeted

342 search and interpretation of OMICS data on male infertility and germ cells in the context of  
343 relevant publications. Moreover, its capacity for aggregation allows for meta-analyses and data  
344 mining with the potential to reveal novel insights into male infertility based on available data.  
345 Ultimately, by combining RGV and MFGA with AI methods, we aim to develop a powerful gene  
346 prioritisation system dedicated to male infertility similar to the GPSy tool (Britto *et al.*, 2012).

347

### 348 **Acknowledgements**

349 We are grateful to all members of the Clinical Research Unit (CRU) 'Male Germ Cells' who  
350 contributed through either identifying relevant publications, supporting publication curation or  
351 providing feedback during the development of the MFGA: Michael Storck, Marius Wöste, Nina  
352 Neuhaus, Sven Berres, Sandra Laurentino, Lina Franziska Lanuza Pérez, Corinna Friedrich,  
353 Maria Schubert, Nikithomas Loges, Isabella Aprea, Jana Emich, Eva Maria Mall, Johanna  
354 Raidt, Nadja Rotte, Alexander Busch, Sara Kim Plutta, Jascha Henseler, Anna Natrup. We  
355 thank Dr. Celeste Brennecka for language editing of the manuscript.

356

### 357 **Authors' roles**

358 H.K. designed and developed the MFGA data model and system architecture, implemented  
359 the tool and drafted the manuscript. J.G. provided critical input during the development of the  
360 MFGA and first drafts of the manuscript. T.D. and F.C. contributed to integrating the mutual  
361 links between RGV and MFGA. M.D. and F.T. contributed to the system design and supervised  
362 the whole study. All authors critically revised the manuscript and approved the final version.

363

### 364 **Funding**

365 This work was carried out within the frame of the German Research Foundation (DFG) Clinical  
366 Research Unit 'Male Germ Cells: from Genes to Function' (CRU326).

367

368 **Conflict of interest**

369 The authors declare no conflicts of interest.

370

371 **References**

- 372 Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, Marshall KA, Phillippy  
373 KH, Sherman PM, Holko M, et al. NCBI GEO: archive for functional genomics data sets--  
374 update. *Nucleic acids research* 2013;**41**:D991-D995.
- 375 Britto R, Sallou O, Collin O, Michaux G, Primig M, Chalmel F. GPSy: a cross-species gene  
376 prioritization system for conserved biological processes--application in male gamete  
377 development. *Nucleic acids research* 2012;**40**:W458-65.
- 378 Budde R, Kautz K, Kuhlenkamp K, Züllighoven H. What is prototyping? *Info Technology &*  
379 *People* 1992;**6**:89-95.
- 380 Canela-Xandri O, Rawlik K, Tenesa A. An atlas of genetic associations in UK Biobank. *Nature*  
381 *genetics* 2018;**50**:1593-1599.
- 382 Chinosi M, Trombetta A. BPMN: An introduction to the standard. *Computer Standards &*  
383 *Interfaces* 2012;**34**:124-134.
- 384 Darde TA, Lecluze E, Lardenois A, Stévant I, Alary N, Tüttelmann F, Collin O, Nef S, Jégou B,  
385 Rolland AD, et al. The ReproGenomics Viewer: a multi-omics and cross-species resource  
386 compatible with single-cell studies for the reproductive science community. *Bioinformatics*  
387 *(Oxford, England)* 2019;**35**:3133-3139.
- 388 Darde TA, Sallou O, Becker E, Evrard B, Monjeaud C, Le Bras Y, Jégou B, Collin O, Rolland  
389 AD, Chalmel F. The ReproGenomics Viewer: an integrative cross-species toolbox for the  
390 reproductive science community. *Nucleic acids research* 2015;**43**:W109-W116.
- 391 Diehl AD, Meehan TF, Bradford YM, Brush MH, Dahdul WM, Dougall DS, He Y, Osumi-  
392 Sutherland D, Ruttenberg A, Sarntivijai S, et al. The Cell Ontology 2016: enhanced content,  
393 modularization, and ontology interoperability. *Journal of biomedical semantics* 2016;**7**:44.
- 394 Gordon VS, Bieman JM. Rapid prototyping: lessons learned. *IEEE Softw.* 1995;**12**:85-95.
- 395 Gremse M, Chang A, Schomburg I, Grote A, Scheer M, Ebeling C, Schomburg D. The  
396 BRENDA Tissue Ontology (BTO): the first all-integrating ontology of all organisms for enzyme  
397 sources. *Nucleic acids research* 2011;**39**:D507-13.
- 398 Hammoud SS, Low DHP, Yi C, Lee CL, Oatley JM, Payne CJ, Carrell DT, Guccione E, Cairns  
399 BR. Transcription and imprinting dynamics in developing postnatal male germline stem cells.  
400 *Genes & development* 2015;**29**:2312-2324.
- 401 HGNC Database. *HGNC Database: retrieved in November 2018*: HUGO Gene Nomenclature  
402 Committee (HGNC), European Molecular Biology Laboratory, European Bioinformatics  
403 Institute (EMBL-EBI), Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, United  
404 Kingdom, 2018.
- 405 Köhler S, Carmody L, Vasilevsky N, Jacobsen JOB, Danis D, Gouridine J-P, Gargano M, Harris  
406 NL, Matentzoglou N, McMurry JA, et al. Expansion of the Human Phenotype Ontology (HPO)  
407 knowledge base and resources. *Nucleic acids research* 2019;**47**:D1018-D1027.
- 408 Lardenois A, Gattiker A, Collin O, Chalmel F, Primig M. GermOnline 4.0 is a genomics gateway  
409 for germline development, meiosis and the mitotic cell cycle. *Database the journal of biological*  
410 *databases and curation* 2010;**2010**:baq030.
- 411 Lonsdale J, Thomas J, Salvatore Mea. The Genotype-Tissue Expression (GTEx) project.  
412 *Nature genetics* 2013;**45**:580-585.



- 413 Maguire M, Bevan N. User Requirements Analysis. In: Hammond J, Gross T, Wesson J (eds).  
414 *Usability*. Boston, MA: Springer US, 2002.
- 415 O'Leary NA, Wright MW, Brister JR, Ciuffo S, Haddad D, McVeigh R, Rajput B, Robbertse B,  
416 Smith-White B, Ako-Adjei D, et al. Reference sequence (RefSeq) database at NCBI: current  
417 status, taxonomic expansion, and functional annotation. *Nucleic acids research* 2016;**44**:D733-  
418 D745.
- 419 Otto M, Thornton J, contributors aB. *Introduction*. Retrieved 09 January 2020. from  
420 <https://getbootstrap.com/docs/4.4/getting-started/introduction/>, 2020.
- 421 Oud MS, Volozonoka L, Smits RM, Vissers LELM, Ramos L, Veltman JA. A systematic review  
422 and standardized clinical validity assessment of male infertility genes. *Human reproduction*  
423 (*Oxford, England*) 2019;**34**:932–941.
- 424 Pivotal Software Inc. *Spring Security*. Retrieved 10 December 2019. from  
425 <https://spring.io/projects/spring-security>, 2019.
- 426 R Development Core Team. *R: A language and environment for statistical computing*. Vienna:  
427 R Foundation for Statistical Computing, 2008.
- 428 Smith B, Ashburner M, Rosse C, Bard J, Bug W, Ceusters W, Goldberg LJ, Eilbeck K, Ireland  
429 A, Mungall CJ, et al. The OBO Foundry: coordinated evolution of ontologies to support  
430 biomedical data integration. *Nat Biotechnol* 2007;**25**:1251–1255.
- 431 Stelzer G, Rosen N, Plaschkes I, Zimmerman S, Twik M, Fishilevich S, Stein TI, Nudel R,  
432 Lieder I, Mazor Y, et al. The GeneCards Suite: From Gene Data Mining to Disease Genome  
433 Sequence Analyses. *Current protocols in bioinformatics* 2016;**54**:1.30.1-1.30.33.
- 434 The Broad Institute of MIT and Harvard. *GTEX Portal*. File: *GTEX\_Analysis\_2016-01-*  
435 *15\_v7\_RNASeQCv1.1.8\_gene\_median\_tpm.gct.gz*. Retrieved 09 December 2019. from  
436 <https://gtexportal.org/home/datasets>, 2019.
- 437 Tüttelmann F, Ruckert C, Röpke A. Disorders of spermatogenesis: Perspectives for novel  
438 genetic diagnostics after 20 years of unchanged routine. *Medizinische Genetik Mitteilungsblatt*  
439 *des Berufsverbandes Medizinische Genetik e.V* 2018;**30**:12–20.
- 440 van der Bijl N, Röpke A, Biswas U, Wöste M, Jessberger R, Kliesch S, Friedrich C, Tüttelmann  
441 F. Mutations in the stromal antigen 3 (STAG3) gene cause male infertility due to meiotic arrest.  
442 *Human reproduction (Oxford, England)* 2019;**34**:2112–2119.
- 443 Yates B, Braschi B, Gray KA, Seal RL, Tweedie S, Bruford EA. Genenames.org: the HGNC  
444 and VGNC resources in 2017. *Nucleic acids research* 2017;**45**:D619-D625.
- 445 Zhang Y, Zhong L, Xu B, Yang Y, Ban R, Zhu J, Cooke HJ, Hao Q, Shi Q.  
446 SpermatogenesisOnline 1.0: a resource for spermatogenesis based on manual literature  
447 curation and genome-wide data mining. *Nucleic acids research* 2013;**41**:D1055-62.
- 448

449 **Table 1. Information classes and items available for each publication.**

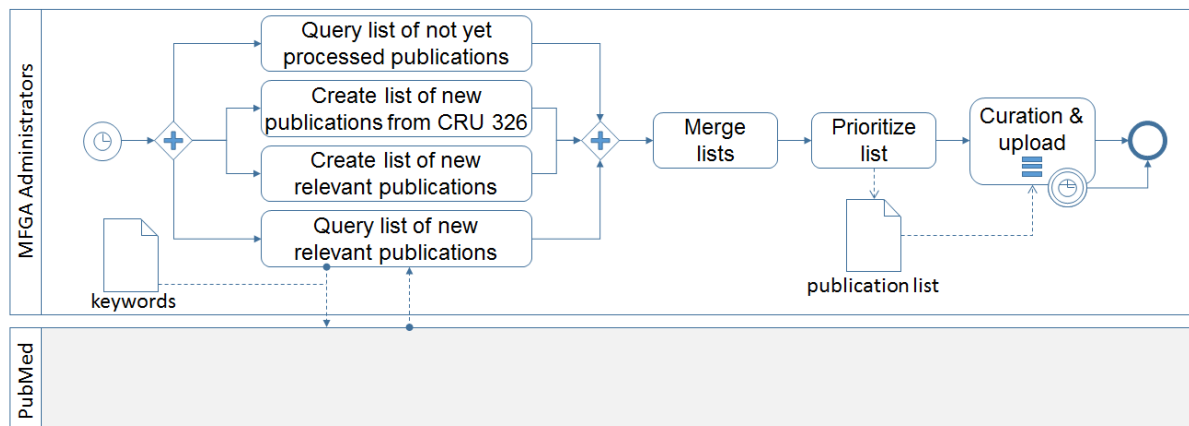
Information Class	Information items <sup>1</sup>	
Publication meta information	title	publishing date
	author	important genes*
	citation	link to publication
	abstract	link to repository*
Data set meta information	title*	OMICS type
	technology*	reference genome*
	data type	species
Processed tissue types	name	maturity*
	description*	potential*
	no of subjects*	species
	no of probes per subject*	reference genome*
Processed cell types	name	maturity*
	description*	potential*
	no of subjects*	species
	no of cells per subject*	reference genome*
Cohort's conditions	Human phenotype ontology* id (Köhler <i>et al.</i> , 2019)	name
	size of cohort*	comment*
Images	title	description*
	annotation of genes or relevant IDs	path
Tables	title	description*
	annotation of standard columns (gene names, relevant IDs, loci, ...)	annotation of all additional columns as numeric or text

450

451 <sup>1</sup>Optional items are marked with \*.

452

453 **Figure 1. The process of identifying, selecting and prioritising novel publications**  
454 **depicted based on Business Process Model and Notation (Chinosi and Trombetta, 2012).**  
455 The process is triggered monthly. Potential publications are collected and merged into a list.  
456 This list is then prioritised by the MFGA team. Subsequently, curation and upload are  
457 performed in the determined order.



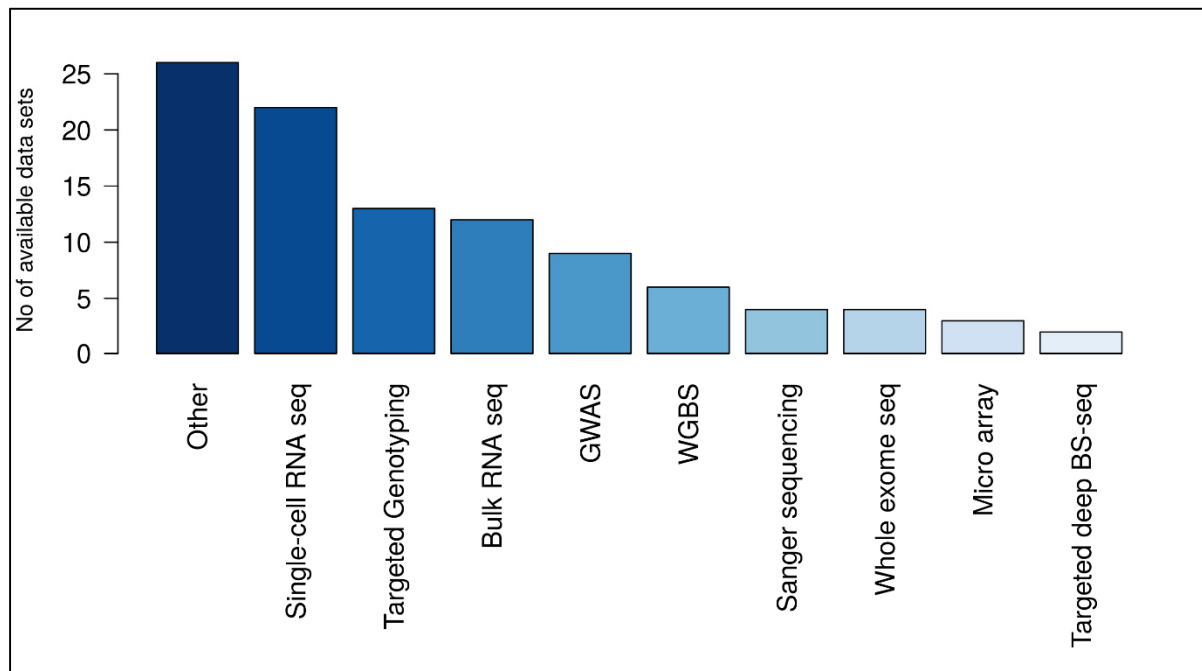
458

459

Male Fertility Gene Atlas (MFGA)

---

460 **Figure 2. Data sets available on the MFGA grouped by technology.** Currently 101 data  
461 sets corresponding to 46 publications are fully curated and publically accessible.  
462 Abbreviations: seq – sequencing, GWAS – genome-wide association study, WGBS – whole  
463 genome bisulfite sequencing, BS – bisulfite. Plot created with R (R Development Core Team,  
464 2008).

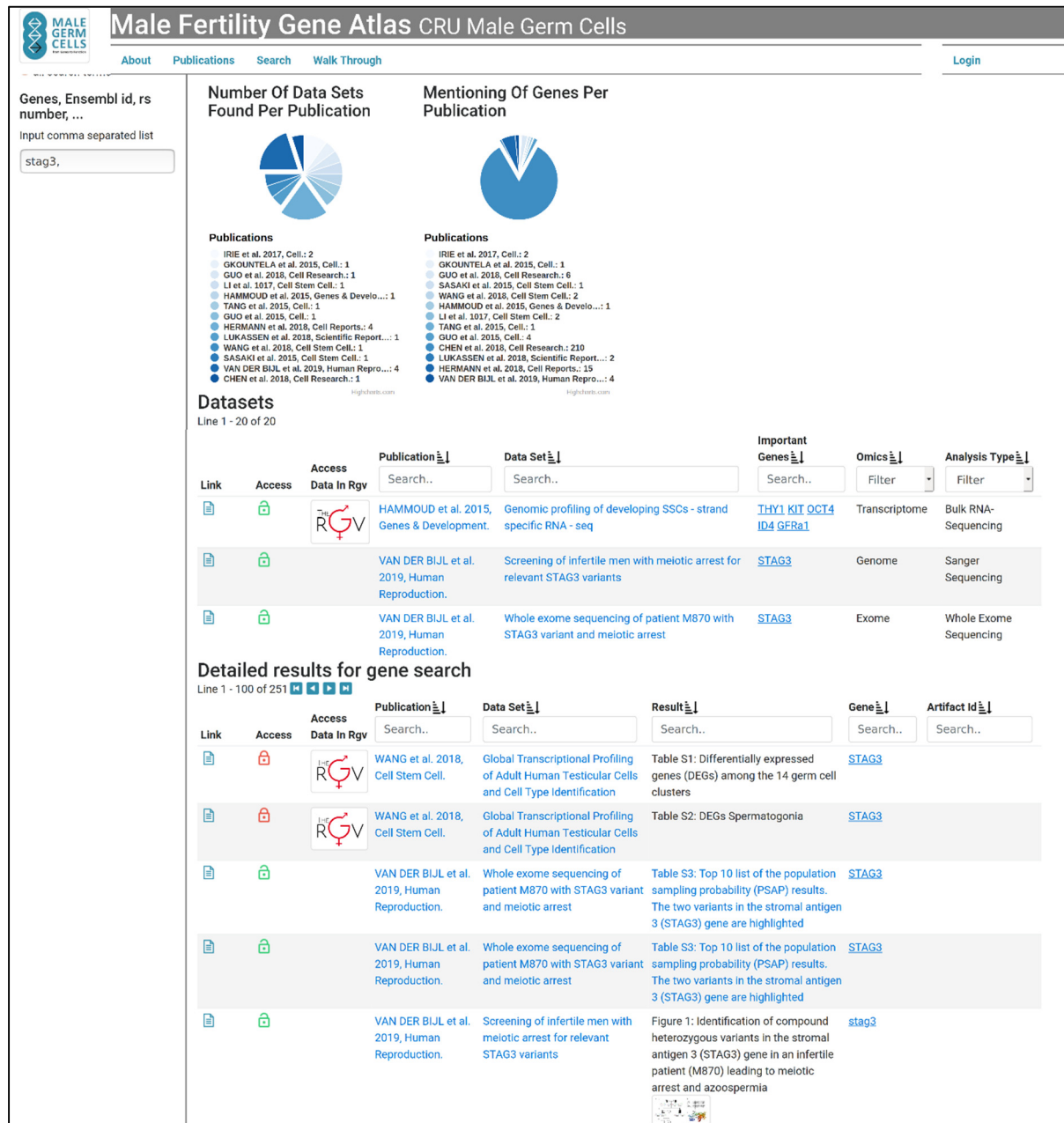


465

466

Male Fertility Gene Atlas (MFGA)

467 **Figure 3. Example screenshot of the MFGA's search functionality** ([https://mfga.uni-](https://mfga.uni-muenster.de/search.html)  
 468 [muenster.de/search.html](https://mfga.uni-muenster.de/search.html)). It shows part of the result of the search for the gene *STAG3* in the  
 469 MFGA database. The side bar on the left contains the search form. Here, *STAG3* is typed into  
 470 the search field and search option 'all search terms' is checked. In the main window, plots  
 471 show the number of returned data sets and the frequency of *STAG3* in data tables of the  
 472 corresponding publications. Below, the returned data sets are presented in a table enriched  
 473 with further meta-information, linking to the overview of data set and publication. The second  
 474 table shows the specific data tables and images in which the gene was found. Here, tables are  
 475 cropped and screenshot is compressed for better representability. See Suppl. Fig. S3 for  
 476 original screenshot. Retrieved 10 January 2020.

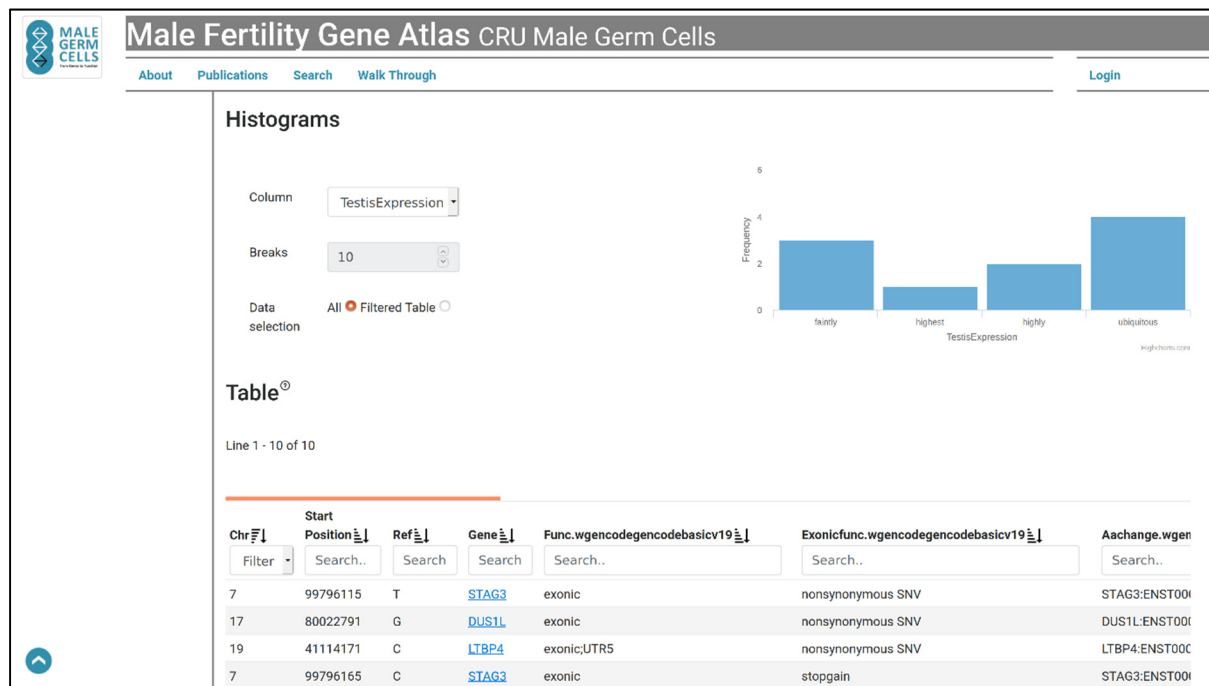


477

478

Male Fertility Gene Atlas (MFGA)

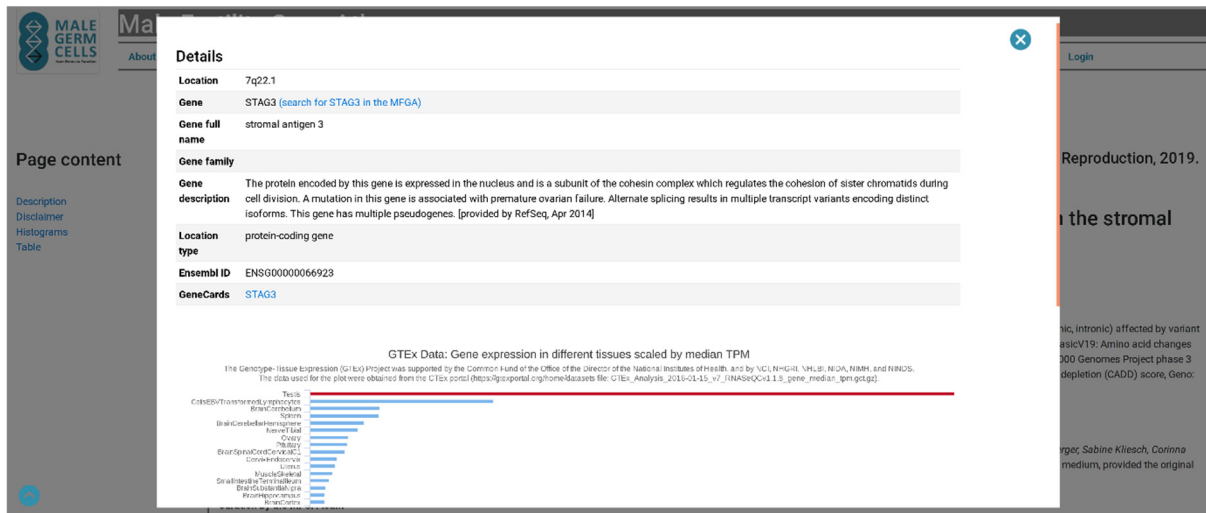
479 **Figure 4. Example screenshot of the MFGA's functionality to represent data from**  
480 **publications** (<https://mfqa.uni-muenster.de/publication/results?ssDetailsId=3968811>). It  
481 shows part of supplementary table 3 of van der Bijl *et al.* (2019) in MFGA format. The table is  
482 represented in an interactive way and can directly be filtered and sorted. Additionally,  
483 histograms can be plotted for the table columns, here, e.g., based on column "Testis  
484 Expression". For the full screenshot see Suppl. Fig. S6. Retrieved 10 January 2020.



485

486

487 **Figure 5. Screenshot of the MFGA showing the overlay for additional information on**  
488 **genes** (<https://mfga.uni-muenster.de/publication/results?ssDetailsId=3968811>). Here, as an  
489 example, the overlay for the gene STAG3 is shown. The textual information originates from  
490 RefSeq (O'Leary *et al.*, 2016) and HGNC (Yates *et al.*, 2017). Additionally, gene expression in  
491 different tissues is shown based on the data from the GTEx project (Lonsdale *et al.*, 2013).  
492 Links for searching the gene in the MFGA database and opening the corresponding  
493 GeneCards entry (Stelzer *et al.*, 2016) are provided. Retrieved 10 January 2020.



494