# Pixel: a content management platform for quantitative omics data

Thomas Denecker[1,*], William Durand[2,*], Julien Maupetit[2,*], Charles Hébert[3], Jean-Michel Camadro[4], Pierre Poulain[4,*,δ] and Gaëlle Lelandais[1,*,δ]

[1] CEA, CNRS, Univ. Paris-Sud, Institute for Integrative Biology of the Cell (I2BC), Gif-sur-Yvette cedex, France

[2] TailorDev SAS, Clermont-Ferrand, France

[3] BIOROSETICS, Houilles, France

[4] CNRS, Univ. Paris Diderot, Institut Jacques Monod (IJM), Paris, France

* These authors contributed equally

δ Corresponding authors: pierre.poulain@univ-paris-diderot.fr ; gaelle.lelandais@u-psud.fr

## ABSTRACT

**Background**. In biology, high-throughput experimental technologies, also referred as "omics" technologies, are increasingly used in research laboratories. Several thousands of gene expression measurements can be obtained in a single experiment. Researchers are routinely facing the challenge to annotate, store, explore and mine all the biological information they have at their disposal. We present here the Pixel web application (Pixel Web App), an original content management platform to help people involved in a multi-omics biological project.

**Methods**. The Pixel Web App is built with open source technologies and hosted on the collaborative development platform GitHub (https://github.com/Candihub/pixel). It is written in Python using the Django framework and stores all the data in a PostgreSQL database. It is developed in the open and licensed under the BSD 3-clause license. The Pixel Web App is also heavily tested with both unit and functional tests, a strong code coverage and continuous integration provided by CircleCI. To ease the development and the deployment of the Pixel Web App, Docker and Docker Compose are used to bundle the application as well as its dependencies.

**Results.** The Pixel Web App offers researchers an intuitive way to annotate, store, explore and mine their multi-omics results. It can be installed on a personal computer or on a server to fit the needs of many users. In addition, anyone can enhance the application to better suit their needs, either by contributing directly on GitHub (encouraged) or by extending Pixel on their own. The Pixel Web App does not provide any computational programs to analyze the data. Still, it helps to rapidly explore and mine existing results and holds a strategic position in the management of research data.

## Introduction

35 
36 In biology, high throughput (HT) experimental technologies - also referred as "omics" - are

37 routinely used in an increasing number of research teams. Financial costs associated to HT

38 experiments have been considerably reduced in the last decade (Hayden, 2014) and the trend

39 in HT sequencing (HTS) is now to acquire benchtop machines designed for individual

40 research laboratories (for instance Illumina NextSeq500 or Oxford Nanopore Technologies

41 MinION, (Blow, 2013)). The number of HT applications in biology has grown so rapidly in

42 the past decade that it is hard to not feel overwhelmed (Hadfield & Retief, 2018)("The data

43 deluge," 2012). It seems possible to address in any organism, any biological question through

44 an "omics" perspective, providing the right HT material and method are found. If HTS is

45 often put at the forefront of "omics" technologies (essentially genomics and

46 transcriptomics, (Reuter, Spacek & Snyder, 2015)), other technologies must be considered.

47 Mass spectrometry (MS) for instance, enables HT identification and quantification of proteins

48 (proteomics). Metabolomics and lipidomics are other derived applications of MS to

49 characterize quantitative changes in small-molecular weight cellular components (Smith et al.,

50 2014). Together, they all account for complementary "omics area" with the advantage to

51 quantify distinct levels of cellular components (transcripts, proteins, metabolites, etc.).

52 Integration of datasets issued from different HT technologies (termed as multi-omics datasets)

53 represents a challenging task from a statistical and methodological point of view (Huang,

54 Chaudhary & Garmire, 2017). It implies the manipulation of two different types of data. The

55 first type is the "primary data", which correspond to raw experimental results. It can be

56 FASTQ files for sequencing technology (Cock et al., 2010) or mzML files for MS (Martens et

57 al., 2011). These files can be stored in public repositories such as SRA (Leinonen et al.,

58 2011), GEO (Clough & Barrett, 2016), PRIDE (Martens et al., 2005) or PeptideAtlas (Desiere

59 et al., 2006). Analyses of primary data rely on standard bioinformatics protocols that for

60 instance, perform quality controls, correct experimental bias or convert files from a specific

61 format to another. A popular tool to analyse primary data is Galaxy (Afgan et al., 2016),

62 which is an open web-based platform. "Secondary data" are produced upon analysis of

63 primary data. It can be the counts of reads per genes for HTS results or the abundance values

64 per proteins for MS results. In multi-omics datasets analysis, combining secondary data is

65 essential to answer specific biological questions. It can be typically, the identification of

66 differentially expressed genes (or proteins) between several cell growth conditions from

67 transcriptomics (or proteomics) datasets, or the identification of cellular functions that are

3

68    over-represented in a list of genes (or proteins). In that respect, secondary data can be

69    analysed and re-analysed within a multitude of analytical strategies, introducing the idea of

70    data analysis cycle. The researcher is thus constantly facing the challenge to

71    properly annotate, store, explore and mine all the biological data he/she has at his/her disposal

72    in a multi-omics project. This challenge is directly related to the ability to extract as much

73    information as possible from the produced data, but also to the crucial question of doing

74    reproducible research.

75    A Nature's survey presented in 2016 indicates that more than 70% of the questioned

76    researchers already experienced an impossibility to reproduce published results, and more

77    than half of them were not able to reproduce their own experiments (Baker, 2016). This last

78    point is intriguing. If experimental biology can be subjected to random fluctuations hardly

79    difficult to control, computational biology should not. Running the same software on the same

80    input data is expected to give the same results. In practice, replication in computational

81    science is harder than people generally think (see (Mesnard & Barba, 2017) as an illustration).

82    It requires to adopt good practices for reproducible-research on a daily basis, and not only

83    when the final results are about to be published. Initiatives to improve computational

84    reproducibility exists (Peng, 2011; Stodden, Guo & Ma, 2013; Vasilevsky et al., 2017;

85    Rougier et al., 2017; Stodden, Seiler & Ma, 2018), and today it is clear that the data alone are

86    not enough to sustain scientific claims. Comments, explanations, software source codes and

87    tests are prerequisites to ensure that an original research can be replicated by anyone, anytime,

88    anywhere.

89    We developed the Pixel web application (Pixel Web App) with these ideas in mind. It is a

90    content management platform to help the researchers involved in a multi-omics biological

91    project, to collaboratively work with their HT data. The Pixel Web App does not store the

92    primary data. It is rather focused on annotation, storage and exploration of secondary data

93    (see **Figure 1**). These explorations represent critical steps to answer biological questions and

94    need to be carefully annotated and recorded to be further exploited in the context of new

95    biological questions. The Pixel Web App helps the researcher to specify necessary

96    information required to replicate multi-omics results. We added an original hierarchical

97    system of tags, which allows to easily explore and select multi-omics results stored in the

98    system and to use them for new interpretations. The Pixel Web App can be installed on any

99    individual computer (for a single researcher for instance), or on a web server for collaborative

100   work between several researchers or research teams. The entire software has been developed

101 with high quality programming standards and complies to major rules of open-source

102 development (Taschuk & Wilson, 2017). The Pixel project is available on GitHub

103 at https://github.com/Candihub/pixel, where full source code and detailed documentation are

104 provided. We present in this article the Pixel Web App design and implementation. We

105 provide a simple case study, emblematic of our daily use of the Pixel Web App, with the

106 exploration of results issued from transcriptomics and proteomics experiments performed in

107 the pathogenic yeast *Candida glabrata*.

## Material and Methods

109 **Stack overview**

110 The Pixel Web App provides researchers an intuitive way to annotate, store, explore and mine

111 their secondary data analyses, in multi-omics biological projects. It is built upon mainstream

112 open source technologies (see **Figure 2**). Source code is hosted on the collaborative

113 development platform GitHub[1] and continuous integration is provided by CircleCI[2]. More

114 precisely, the Pixel Web App uses the Python Django framework. This framework is based on

115 a model-template-view architecture pattern, and data are stored in a PostgreSQL[3] database.

116 We have built a docker image for the Pixel Web App. Other containers, Nginx (to serve the

117 Django application) and PostgreSQL rely on official docker images. Each installation /

118 deployment will result in the creation / execution of three docker instances: one for the Pixel

119 Web App, one for the PostgreSQL database and one for the Nginx web server. In case of

120 multiple installations, each trio of docker instances is fully isolated, meaning that data are not

121 shared across multiple Pixel Web App installations.

122 **Technical considerations**

123 • Docker images

124 The Pixel Web App is built on containerization paradigm (see **Figure 2**). It relies

125 on Docker[4], *i.e.* a tool which packages an application and its dependencies in an image that

126 will be run as a container. Docker helps developers to build self-contained images to run a

127 software. These images are downloaded on the host system and used to build the Pixel Web

128 App.

---

[1] https://github.com/
[2] https://circleci.com/
[3] https://www.postgresql.org/
[4] https://www.docker.com/

129      •    Minimal configuration and dependencies

130   The Pixel Web App can be deployed on Linux and MacOS operating systems (OS).

131   Deployment on Windows is possible, but this situation will not be described here. Minimal

132   requirements are: *(i)* 64 bits Unix-based OS (Linux / MacOS), *(ii)* Docker community edition

133   > v18, *(iii)* Internet access (required in order to download the Docker images) and *(iv)*

134   [optional] a web server (Apache or Nginx) configured as a reverse proxy.

135   **Installation**

136   A step-by-step tutorial to deploy the Pixel Web App can be found in the project repository[5]

137   together with a deploy script. To summarize, this script runs the following steps:

138      ➢   Pull a tagged image of Pixel (web, see docker-composer file),

139      ➢   Start all instances (web, db and proxy) recreating the proxy and web instances. Collect

140         all static files from the Django app. These files will be served by the proxy instance.

141      ➢   Migrate the database schema if needed (to preserve existing data).

142   Note that further technical considerations and full documentation can be found on GitHub

143   repository associated to the Pixel project[6].

144   # Results

145   **Definition of terms: Omics Unit, Pixel and Pixel Set**

146   In the Pixel Web App, the term "Omics Unit" refers to any cellular component, from any

147   organism, which is of interest for the user. The type of Omics Unit depends on the HT

148   experimental technology (transcriptomic, proteomic, metabolomic, etc.) from which primary

149   and secondary datasets were collected and derived (**Figure 1**A). In this context, classical

150   Omics Units can be transcripts or proteins, but any other cellular component can be defined

151   as, for instance, genomic regions with "peaks" in case of ChIPseq data analyses (Merhej et al.,

152   2014). A "Pixel" refers to a quantitative measurement of a cellular activity associated to a

153   single Omics Unit, together with a quality score (see **Figure 1**A). Quantitative measurement

154   and quality score are results of statistical analyses performed on secondary

155   datasets, *e.g.* search for differentially expressed genes (Seyednasrollah, Laiho & Elo, 2015). A

156   set of Pixels obtained from a single secondary data analysis of HT experimental results is

157   referred as a "Pixel Set" (see **Figure 1**A). Pixel Sets represent the central information in the

---

[5] https://github.com/Candihub/pixel/blob/master/docs-install/how-to-install.md
[6] https://github.com/Candihub/pixel/tree/master/docs

158  Pixel Web App and functionalities to annotate, store, explore and mine multi-omics biological

159  data were designed according to this concept (see below).

160  **Functionalities to annotate, store, explore and mine Pixel Sets**

161  Pixel Sets are obtained from secondary data analyses (see **Figure 1**A). Their manipulation

162  with the Pixel Web App consists in *(i)* their annotation, *(ii)* their storage in a database, *(iii)*

163  their exploration and *(iv)* their mining (see **Figure 1**C). This represents a cycle of multiple

164  data analyses, which is essential in any multi-omics biological project. These different steps

165  are detailed in the following.

166  • Annotation of Pixel Sets

167  Annotation of Pixel Sets consists in tracking important details of Pixel Set production. For

168  that, Pixel Sets are associated with metadata, *i.e.* supplementary information linked to the

169  Pixel Sets. We defined minimal information necessary for relevant annotations of Pixel Sets

170  (see **Figure 3**). "Species", "Strain", "Omics Unit Type" and "Omics Area" are mandatory

171  information that must be specified *before* a new Pixel Set submission (highlighted in blue,

172  **Figure 3**). They refer to general information related to the multi-omics biological project on

173  which the researcher is working on: *(i)* the studied organism and its genetic background

174  (Species and Strain, *e.g. Candida glabrata* and ATCC2001), *(ii)* the type of monitored

175  cellular components (Omics Unit Type, *e.g.* mRNA, protein) and *(iii)* the nature of the

176  experimental HT technology (Omics Area, *e.g.* RNA sequencing, mass spectrometry). All

177  Omics Units must be declared in the Pixel Web App before new Pixel Set submission. They

178  must be defined with a short description and a link to a reference database. "Experiment" and

179  "Analysis" are Pixel Set mandatory information, input during the submission of new Pixel

180  Sets in the Pixel Web App (highlighted in orange, **Figure 3**). They include respectively the

181  detailed description of the experimental strategy that was applied to generate primary and

182  secondary data sets (Experiment) and the detailed description of the computational procedures

183  that were applied to obtain Pixel Sets from secondary data set (Analysis). Information

184  regarding the researcher who performed the analyses is referred as "Pixeler".

185  • Storage of Pixel Sets in the database

186  Import of new Pixel Sets in the Pixel Web App requires the user to follow a workflow for data

187  submission.  It corresponds to six successive steps that are explained below (**Figure 4**A).

188  1. The "Download" step consists in downloading a template Excel file from the Pixel

189  Web App (see **Figure 4**B). In this file, multiple-choice selections are proposed for

7

190    "Species", "Strain", "Omics Unit Type" and "Omics Area" fields. These choices

191    reflect what is currently available in the database and can be easily expanded. User

192    must fill other annotation fields related to the "Experiment", "Analysis" and "Pixeler"

193    information. The Excel file is next bundled into a ZIP archive with the secondary data

194    file (in tab-separated values format), the user notebook (R markdown[7] or Jupyter

195    notebook[8] for instance) that contains the code used to produce the Pixel Sets from the

196    secondary data file.

197    2.  The "Upload" step consists in uploading the ZIP file in the Pixel Web App.

198    3.  The step "Meta" consists in running an automatic check of the imported file

199    integrity (md5sum checks are performed, Excel file version is verified, etc.). Note that

200    no information is imported in the database at this stage, but a careful inspection of

201    all Omics Units listed in the submitted Pixel Sets is done. This is why Omics Units

202    need to be pre-registered in the Pixel Web App (see previous section).

203    4.  In "Annotation" step, the annotations of Pixel Sets found in the Excel file (see **Figure**

204    **4**C) are controlled and validated by the user.

205    5.  Next, the "Tags" step is optional. It gives the opportunity to the user to add tags to the

206    new Pixel Sets (see **Figure 4**C), that could be helpful for further Pixel Set explorations

207    (see next section).

208    6.  The final step "Import archive" consists in importing all Pixel Sets in the database,

209    together with annotations and tags.

210    Note that the procedure of importing meta data as an Excel file has been inspired from the

211    import procedure widely used in GEO (Clough & Barrett, 2016).

212    •   Exploration of Pixel Sets

213    The Pixel Web App aims to help researchers to mine and integrate multiple Pixel Sets stored

214    in the system. We developed a dedicated web interface to explore all the Pixel Sets stored in a

215    particular Pixel instance (see **Figure 5**). The upper part named "Selection" lists a group of

216    Pixel Sets selected by the user for further explorations (**Figure 5**A). The middle part named

217    "Filters" lists the Pixel database contents regarding the Species, Omics Unit Types, Omics

218    Areas and Tags annotation fields. The user can select information (*Candida glabrata*

219    and modified pH here), search and filter the Pixel Sets stored in the database (**Figure 5**B).

220    The lower part is a more flexible search field in which keywords can be type. These keywords

---

[7] https://rmarkdown.rstudio.com/
[8] http://jupyter.org/

221    are searched in the Analysis and Experiment detailed description fields as illustrated here

222    with LIMMA. The web interface also comprised detailed information for the selected subset

223    of Pixel Sets with for instance, distributions of values and quality scores and a list of

224    individual Omics Unit shown at the bottom of the page (**Figure 5**C). Note that tags have been

225    implemented to offer to the user a versatile yet robust annotation of Pixel Sets. They are

226    defined during the import process, but they can be modified at any time through the Pixel web

227    interface. Once searched, matching Pixel Sets are gathered in a table that can be exported.

228    **A case study in the pathogenic yeast _Candida glabrata_**

229    The yeast _Candida glabrata_ (_C. glabrata_) is a fungal pathogen of human (Bolotin-Fukuhara

230    & Fairhead, 2014). It has been reported as the second most frequent cause of invasive

231    infections due to Candida species, _i.e._ candidemia, arising especially in patients with

232    compromised immunity (HIV virus infection, cancer treatment, organ transplantation, etc.).

233    Candidemia remains a major cause of morbidity and mortality in the healthcare

234    structures (Horn et al., 2009; Pfaller et al., 2012). The genome of _Candida glabrata_ has been

235    published in 2004 (Dujon et al., 2004). Its size is 12.3 Mb with 13 chromosomes and is

236    composed of ~5200 coding regions.  Our research team is familiar with functional genomic

237    studies in _C. glabrata_. In collaboration with experimental biologists, we published in the past

238    ten years half dozen of articles, in which HT technologies were used (Lelandais et al., 2008;

239    Goudot et al., 2011; Merhej et al., 2015, 2016; Thiébaut et al., 2017). In our lab, the Pixel

240    Web App is installed locally and store all the necessary genomics annotations to manage any

241    multi-omics datasets in this species.

242    As a case study, we decided to present how the Pixel Web App can be helpful to answer a

243    specific biological question with only a few mouse clicks. As a biological question, we

244    wanted to identify the genes in the entire _C. glabrata_ genome: _(i)_ which are annotated as

245    involved in the yeast pathogenicity and _(ii)_ for which the expression is significantly modified

246    in response to an environmental stress induced by alkaline pH. Indeed, during a human host

247    infection, _C. glabrata_ has to face important pH fluctuations (see (Ullah et al., 2013; Brunke &

248    Hube, 2013; Linde et al., 2015) for more detailed information). Understanding the molecular

249    processes that allow the pathogenic yeast _C. glabrata_ to adapt extreme pH situations is

250    therefore of medical interest to better understand host-pathogen interaction (Linde et al.,

251    2015).

9

252   In a paper published in 2015, Linde *et al.* provided a detailed RNAseq based analysis of the

253   transcriptional landscape of *C. glabrata* in several growth conditions, including pH shift

254   experiments (Linde et al., 2015).  The primary dataset (RNAseq fastq files) is available in the

255   Gene Expression Omnibus (Clough & Barrett, 2016) under accession number GSE61606. The

256   secondary dataset (log2 Fold Change values) is available in Supplementary Table S1 on the

257   journal website[9]. A first Pixel Set (labelled A) was created from this secondary dataset,

258   annotated and imported into our Pixel Web App instance, following the procedure previously

259   described. The associated ZIP archive is provided as supplemental file, along with the all the

260   details related to the experiment set up and the analysis. The Pixel Set A thus illustrates how

261   publicly available data can be managed with the Pixel Web App. In our laboratory, we

262   performed mass spectrometry experiments that also include pH shift (unpublished results, but

263   ZIP archive of the data is provided as supplemental file). Secondary dataset issued from these

264   experiments leads to the Pixel Set B. Pixel Sets A and B comprise 5,253 Pixels and 1,879

265   Pixels (**Figure 6**).

266   Transcriptomics (Pixel Set A) and proteomics (Pixel Set B) are interesting complementary

267   multi-omics information that can be easily associated and compared with the Pixel Web App.

268   In that respect, tags allowed to rapidly retrieve them using the web interface, applying the

269   keywords "Candida glabrata" and "alkaline pH" (**Figure 6**, Step 1). As we wanted to limit the

270   analysis to the *C. glabrata* genes potentially involved in the yeast pathogenesis, a filter could

271   be used to only retain the Omics Units for which the keyword "pathogenicity" is written in

272   their description field (see **Figure 6**, Step 2). As a result, a few numbers of Pixels were thus

273   selected, respectively 17 in Pixel Set A and 6 in Pixel Set B. The last step consists in

274   combining the mRNA and protein information (see **Figure 6**, Step 3). For that a table

275   comprising the multi-pixel sets can be automatically generated and easily exported. We

276   present **Table 1** five genes for which logFC values were obtained both at the mRNA and the

277   protein levels, and for which statistical p-values were significant ($< 0.05$). Notably two genes

278   (CAGL0I02970g and CAGL0L08448g, lines 3 and 5 in **Table 1**) exhibited opposite logFC

279   values, *i.e.* induction was observed at the mRNA level whereas repression was observed at the

280   protein levels. Such observations can arise from post-translational regulation processes or

281   from possible experimental noise, which could explain approximative mRNA or protein

282   quantifications. In both cases, further experimental investigations are required. The three

283   other genes (CAGL0F04807g, CAGL0F06457g and CAGL0I10516g, underlined in grey

---

[9] https://academic.oup.com/nar/article/43/3/1392/2411170

10

284    **Table 1**) exhibited multi-omics coherent results and significant inductions were observed at

285    the mRNA and protein levels. Again, further experimental investigations are required to fully

286    validated these observations. Still, it is worth noting that the gene CAGL0F04807g, is

287    described as "uncharacterized" in the Candida Genome Database [10]. Considering that logFC

288    values for this gene are particularly high ($> 1$), such an observation represents a good starting

289    point to refine the functional annotation of this gene, clearly supporting the hypothesis that is

290    has a role in the ability of *C. glabrata* to deal with varying pH situations.

**Software Availability**

292    Pixel is released under the open-source 3-Clause BSD license

293    (https://opensource.org/licenses/BSD-3-Clause). Its source code can be freely downloaded

294    from the GitHub repository of the project: https://github.com/Candihub/pixel. In addition, the

295    present version of Pixel (4.0.4) is also archived in the digital repository Zenodo

296    (https://doi.org/10.5281/zenodo.1434316).

## Discussion

298    In this article, we introduced the principle and the main functionalities of the Pixel Web App.

299    With this application, our aim was to develop a tool to support on a daily basis, the biological

300    data mining in our multi-omics research projects. It is our experience that research studies in

301    which HT experimental strategies are applied, require much more time to analyse and

302    interpret the data, than to experimentally generate the data. Testing multiple bioinformatics

303    tools and statistical approaches is a critical step to fully understand the meaning of a

304    biological dataset and in this context, the annotation, the storage and the ability to easily

305    explore the all results obtained in a laboratory can be the decisive steps to the success of the

306    entire multi-omics project.

307    The data modelling around which the Pixel Web App was developed, has been conceived to

308    find a compromise between a too detailed and precise description of the data (which could

309    discourage the researchers of systematically use the application after each of their analyses)

310    and a too short and approximate description of the data (which could prevent the

311    perfect reproduction of the results by anyone). Also, an attention has been paid to allow

312    heterogeneous data, *i.e.* different Omics Unit Type quantified in different Omics Area, to be

313    stored in a coherent and flexible way. The Pixel Web App does not provide any

314    computational programs to analyse the data. Still, it allows to explore existing results in a

---

[10] http://www.candidagenome.org/cgi-bin/locus.pl?locus=CAGL0F04807g&organism=C_glabrata_CBS138

315    laboratory and to rapidly combine them for further investigations (using for instance the

316    Galaxy platform or any other data analysis tool).

317    Therefore, the Pixel Web App holds a strategic position in the data management in a research

318    laboratory, *i.e.* as the starting point but also at the final point of all new data explorations. It

319    also helps data analysis reproducibility and gives a constant feedback regarding the frequency

320    of the data analysis cycles; the nature of the import and export data sets as well as full

321    associated annotations. It is thus expected that the content of different Pixel Web App

322    instance will evolve with time, according to the type of information stored in the system and

323    the scientific interests of a research team.

## Conclusion

325    The Pixel Web App is freely available to any interested people. The initial installation on a

326    personal workstation required IT support from a bioinformatician, but once this is done, all

327    administration tasks can be performed through the Web Interface. This is of interest for user

328    with a few technical skills. We chose to work exclusively with open source technologies and

329    our GitHub repository is publicly accessible[11]. We thus hope that the overall quality of the

330    Pixel Web App source code and documentation will be guaranteed over time, through the

331    shared contributions of other developers.

## Figure and table legends

333    **Figure 1: Dataset flow through the Pixel Web App**. (A) Different types of datasets, which are
334    managed in a multi-omics biological project. Primary and secondary datasets are two types of
335    information arising from HT experimental technologies (see the section **Introduction**). Only
336    secondary data and their associated Pixel Sets are stored in the Pixel Web App. Note that several Pixel
337    Sets can emerge from multiple secondary data analyses. They comprise quantitative values (Value)
338    together with quality scores (QS) for several hundred of different "Omics Units" elements (for
339    instance mRNA or proteins, see the main text). Omics Units are identified with a unique identifier
340    (ID). (B) Screenshot of the home page of the Pixel web interface. (C) Schematic representation of the
341    data analysis cycles that surrounds the integration of Pixel Sets in the Pixel Web App (see the main
342    text).

343    **Figure 2: Stack overview of the Pixel Web App**. Open source solutions used to develop Pixel are
344    shown here. They are respectively used for the software development and test (blue section), the data
345    storage (green section) and the web application for both staging and production (orange section).

346    **Figure 3: Data modelling in the Pixel Web App**. The Pixel Set is the central information (see **Figure
347    1**A), the corresponding table in the model is highlighted in red. Information that is required *before*
348    Pixel Set import in the Pixel Web App is surrounded in blue, whereas information required *during*
349    Pixel Set import is highlighted in orange. Other tables are automatically updated during the Pixel Web

---

[11] https://github.com/Candihub/pixel

350     App data analysis life cycle (see **Figure 1**C). Enlarge version of this picture together with full
351     documentation is available online[12].

352     **Figure 4: Procedure to import new Pixel Sets in the Pixel Web App**. (A) New data-sets are
353     submitted following a dedicated workflow that comprised 6 successive actions named "Download",
354     "Upload", "Meta", "Validation", "Tags" and "Import archive" (see 1). Several files are required (see
355     2): the secondary data from which the Pixel Sets were calculated, the notebook in which the procedure
356     to compute Pixel Sets from secondary data is described and the Pixel Set files (2 files in this example).
357     A progression bar allows the user to follow the sequence of the submission process. (B) Excel
358     spreadsheet in which annotations of Pixel Sets are written. Information related to the Experiment (see
359     1), the Analysis (see 2) and the Pixel datasets (see 3) is required. Note that this file must be
360     downloaded at the first step of the submission process ("Download", see A), allowing several cells to
361     be pre-filled with annotations stored in the database (see 4 as an illustration, with Omics area
362     information). (C) All information filled in the Excel file (see B) is extracted and can be modified
363     anytime through a dedicated web page as shown here. User can edit the Pixel Set (see 1), edit the
364     analysis (see 2), edit the experiment (see 3) and add "Tags" (see 4). The Tags are of interest to further
365     explore Pixel Sets in the Pixel Web App.

366     **Figure 5 : Functionalities to explore the Pixel Sets stored in the Pixel Web App**. (A) Screenshot of
367     the exploration menu available *via* the web interface. (B) Screenshot of the table that comprises all
368     Pixel Sets, which match the filter criteria (see A). Particular Pixel Sets can be selected here (for
369     instance "Pixel_C10.txt" and "Pixel_C60.txt"). They will therefore appear in the "Selection" list (see
370     A). (C) Screenshot of the web interface that gives detailed information for the selected subset of Pixel
371     Sets (see A). Distribution of values and quality scores are shown and individual Omics Unit are listed
372     at the bottom of the page.

373     **Figure 6: Case study in the pathogenic yeast *Candida glabrata***. Our Pixel Web App was explored
374     with the keywords "Candida glabrata" and "alkaline pH". Two Pixel Sets were thus identified because
375     of their tags. Two other tags were identical between the two Pixel Sets ("WT" and "logFC"), indicating
376     that *(i) C. glabrata* strains are the same, *i.e.* Wild Type, and *(ii)* Pixel values are of the same
377     type, *i.e.* log Fold Change. Notably Pixel Set A is based on transcriptomics experiments (RNAseq, see
378     the main text), whereas Pixel Set B is based on proteomics experiments (mass spectrometry, see the
379     main text). Omics Unit were next explored searching the keyword "pathogenesis" in their description
380     fields (coming from the CGD database (Skrzypek et al., 2017)). This results in the identification of 17
381     Pixels (respectively 6 Pixels) in transcriptomics (respectively proteomics) results. They were
382     combined and exported from the Pixel Web App, hence starting a new data analysis cycle.

383     **Table 1:  Detailed information regarding the Omics Unit identified in the *C. glabrata* case**
384     **study**. The two first column give Omics Unit information as described in the Candida Genome
385     Database (Skrzypek et al., 2017). All the description fields comprise the keyword "pathogenesis" (in
386     bold). LogFC values measured in transcriptomic (Pixel Set A) and proteomic (Pixel Set B)
387     experiments are shown in the third and fourth columns. Quality scores (QS) are following logFC
388     values. They are p-values coming from the differential analysis of logFC replicates. The entire table of
389     multi-pixel sets is available in supplementary data.

# References

391     Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Čech M, Chilton J, Clements D,
392        Coraor N, Eberhard C, Grüning B, Guerler A, Hillman-Jackson J, Von Kuster G, Rasche E,
393        Soranzo N, Turaga N, Taylor J, Nekrutenko A, Goecks J. 2016. The Galaxy platform for
394        accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic acids*
395        *research* 44:W3–W10. DOI: 10.1093/nar/gkw343.

396     Baker M. 2016. 1,500 scientists lift the lid on reproducibility. *Nature* 533:452–4. DOI:
397        10.1038/533452a.

---

[12] https://github.com/Candihub/pixel/blob/master/docs/pixel-db.pdf

398    Blow N. 2013. A sequencer in every lab. *BioTechniques* 55:284. DOI: 10.2144/000114107.

399    Bolotin-Fukuhara M, Fairhead C. 2014. Candida glabrata: a deadly companion? *Yeast (Chichester,*
400        *England)* 31:279–88. DOI: 10.1002/yea.3019.

401    Brunke S, Hube B. 2013. Two unlike cousins: *Candida albicans* and *C. glabrata* infection strategies.
402        *Cellular Microbiology* 15:701–708. DOI: 10.1111/cmi.12091.

403    Clough E, Barrett T. 2016. The Gene Expression Omnibus Database. *Methods in molecular biology*
404        *(Clifton, N.J.)* 1418:93–110. DOI: 10.1007/978-1-4939-3578-9_5.

405    Cock PJ, Fields CJ, Goto N, Heuer ML, Rice PM. 2010. The Sanger FASTQ file format for sequences
406        with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic acids research* 38:1767–
407        71. DOI: 10.1093/nar/gkp1137.

408    Desiere F, Deutsch EW, King NL, Nesvizhskii AI, Mallick P, Eng J, Chen S, Eddes J, Loevenich SN,
409        Aebersold R. 2006. The PeptideAtlas project. *Nucleic Acids Research* 34:D655–D658. DOI:
410        10.1093/nar/gkj040.

411    Dujon B, Sherman D, Fischer G, Durrens P, Casaregola S, Lafontaine I, De Montigny J, Marck C,
412        Neuvéglise C, Talla E, Goffard N, Frangeul L, Aigle M, Anthouard V, Babour A, Barbe V,
413        Barnay S, Blanchin S, Beckerich J-M, Beyne E, Bleykasten C, Boisramé A, Boyer J, Cattolico L,
414        Confanioleri F, De Daruvar A, Despons L, Fabre E, Fairhead C, Ferry-Dumazet H, Groppi A,
415        Hantraye F, Hennequin C, Jauniaux N, Joyet P, Kachouri R, Kerrest A, Koszul R, Lemaire M,
416        Lesur I, Ma L, Muller H, Nicaud J-M, Nikolski M, Oztas S, Ozier-Kalogeropoulos O, Pellenz S,
417        Potier S, Richard G-F, Straub M-L, Suleau A, Swennen D, Tekaia F, Wésolowski-Louvel M,
418        Westhof E, Wirth B, Zeniou-Meyer M, Zivanovic I, Bolotin-Fukuhara M, Thierry A, Bouchier
419        C, Caudron B, Scarpelli C, Gaillardin C, Weissenbach J, Wincker P, Souciet J-L. 2004. Genome
420        evolution in yeasts. *Nature* 430:35–44. DOI: 10.1038/nature02579.

421    Goudot C, Etchebest C, Devaux F, Lelandais G. 2011. The Reconstruction of Condition-Specific
422        Transcriptional Modules Provides New Insights in the Evolution of Yeast AP-1 Proteins. *PloS*
423        *one* 6:e20924. DOI: 10.1371/journal.pone.0020924.

424    Hadfield J, Retief J. 2018. A profusion of confusion in NGS methods naming. *Nature Methods* 15:7–8.
425        DOI: 10.1038/nmeth.4558.

426    Hayden EC. 2014. The $1,000 genome. *Nature* 507:294. DOI: 10.1038/507294a.

427    Horn DLL, Neofytos D, Anaissie EJJ, Fishman JAA, Steinbach WJJ, Olyaei AJJ, Marr KAA, Pfaller
428        MAA, Chang CC-H, Webster KMM. 2009. Epidemiology and Outcomes of Candidemia in 2019
429        Patients: Data from the Prospective Antifungal Therapy Alliance Registry. *Clinical Infectious*
430        *Diseases* 48:1695–1703. DOI: 10.1086/599039.

431    Huang S, Chaudhary K, Garmire LX. 2017. More is better: Recent progress in multi-omics data
432        integration methods. *Frontiers in Genetics* 8:1–12. DOI: 10.3389/fgene.2017.00084.

433    Leinonen R, Sugawara H, Shumway M, International Nucleotide Sequence Database Collaboration.
434        2011. The sequence read archive. *Nucleic acids research* 39:D19-21. DOI: 10.1093/nar/gkq1019.

435    Lelandais G, Tanty V, Geneix C, Etchebest C, Jacq C, Devaux F. 2008. Genome adaptation to
436        chemical stress: clues from comparative transcriptomics in Saccharomyces cerevisiae and
437        Candida glabrata. *Genome biology* 9:R164. DOI: 10.1186/gb-2008-9-11-r164.

438    Linde JJ, Duggan SS, Weber M, Horn F, Sieber P, Hellwig D, Riege K, Marz M, Martin R, Guthke R,
439        Kurzai O. 2015. Defining the transcriptomic landscape of Candida glabrata by RNA-Seq.
440        *Nucleic Acids Research* 43:1392–1406. DOI: 10.1093/nar/gku1357.

441 Martens L, Chambers M, Sturm M, Kessner D, Levander F, Shofstahl J, Tang WH, Römpp A,
442     Neumann S, Pizarro AD, Montecchi-Palazzi L, Tasman N, Coleman M, Reisinger F, Souda P,
443     Hermjakob H, Binz P-A, Deutsch EW. 2011. mzML—a Community Standard for Mass
444     Spectrometry Data. *Molecular & Cellular Proteomics* 10:R110.000133. DOI:
445     10.1074/mcp.R110.000133.

446 Martens L, Hermjakob H, Jones P, Adamski M, Taylor C, States D, Gevaert K, Vandekerckhove J,
447     Apweiler R. 2005. PRIDE: The proteomics identifications database. *PROTEOMICS* 5:3537–
448     3545. DOI: 10.1002/pmic.200401303.

449 Merhej J, Delaveau T, Guitard J, Palancade B, Hennequin C, Garcia M, Lelandais G, Devaux F. 2015.
450     Yap7 is a Transcriptional Repressor of Nitric Oxide Oxidase in Yeasts, which arose from
451     Neofunctionalization after Whole Genome Duplication. *Molecular Microbiology* 96:n/a-n/a.
452     DOI: 10.1111/mmi.12983.

453 Merhej J, Frigo A, Crom S Le, Camadro JJ, Le Crom S, Camadro JJ, Devaux F, Lelandais G. 2014.
454     bPeaks□: a bioinformatics tool to detect transcription factor binding sites from ChIPseq data in
455     yeasts and other organisms with small genomes. *Yeast* 31:375–391. DOI: 10.1002/yea.

456 Merhej J, Thiebaut A, Blugeon C, Pouch J, Ali Chaouche MEA, Camadro J-M, Le Crom S, Lelandais
457     G, Devaux F. 2016. A Network of Paralogous Stress Response Transcription Factors in the
458     Human Pathogen Candida glabrata. *Frontiers in Microbiology* 7:1–16. DOI:
459     10.3389/fmicb.2016.00645.

460 Mesnard O, Barba LA. 2017. Reproducible and Replicable Computational Fluid Dynamics: It's
461     Harder Than You Think. *Computing in Science & Engineering* 19:44–55. DOI:
462     10.1109/MCSE.2017.3151254.

463 Peng RD. 2011. Reproducible research in computational science. *Science (New York, N.Y.)* 334:1226–
464     7. DOI: 10.1126/science.1213847.

465 Pfaller M, Neofytos D, Diekema D, Azie N, Meier-Kriesche HU, Quan SP, Horn D. 2012.
466     Epidemiology and outcomes of candidemia in 3648 patients: Data from the Prospective
467     Antifungal Therapy (PATH Alliance) registry, 2004-2008. *Diagnostic Microbiology and
468     Infectious Disease* 74:323–331. DOI: 10.1016/j.diagmicrobio.2012.10.003.

469 Reuter JAA, Spacek D V., Snyder MPP. 2015. High-Throughput Sequencing Technologies. *Molecular
470     Cell* 58:586–597. DOI: 10.1016/j.molcel.2015.05.004.

471 Rougier NP, Hinsen K, Alexandre F, Arildsen T, Barba LA, Benureau FCYY, Brown CT, de Buyl P,
472     Caglayan O, Davison AP, Delsuc M-AA, Detorakis G, Diem AK, Drix D, Enel P, Girard B,
473     Guest O, Hall MG, Henriques RN, Hinaut X, Jaron KS, Khamassi M, Klein A, Manninen T,
474     Marchesi P, McGlinn D, Metzner C, Petchey OL, Plesser HE, Poisot T, Ram K, Ram Y, Roesch
475     E, Rossant C, Rostami V, Shifman A, Stachelek J, Stimberg M, Stollmeier F, Vaggi F, Viejo G,
476     Vitay J, Vostinar AE, Yurchak R, Zito T. 2017. Sustainable computational science: the
477     ReScience initiative. *PeerJ Computer Science* 3:1–8. DOI: 10.7717/peerj-cs.142.

478 Seyednasrollah F, Laiho A, Elo LL. 2015. Comparison of software packages for detecting differential
479     expression in RNA-seq studies. *Briefings in Bioinformatics* 16:59–70. DOI: 10.1093/bib/bbt086.

480 Skrzypek MS, Binkley J, Binkley G, Miyasato SR, Simison M, Sherlock G. 2017. The Candida
481     Genome Database (CGD): Incorporation of Assembly 22, systematic identifiers and visualization
482     of high throughput sequencing data. *Nucleic Acids Research* 45:D592–D596. DOI:
483     10.1093/nar/gkw924.

484 Smith R, Mathis A, Ventura D, Prince J. 2014. Proteomics, lipidomics, metabolomics: A mass
485     spectrometry tutorial from a computer scientist's point of view. *BMC Bioinformatics* 15:S9. DOI:

15

486    10.1186/1471-2105-15-S7-S9.

487    Stodden V, Guo P, Ma Z. 2013. Toward Reproducible Computational Research: An Empirical
488        Analysis of Data and Code Policy Adoption by Journals. *PloS one* 8:e67111. DOI:
489        10.1371/journal.pone.0067111.

490    Stodden V, Seiler J, Ma Z. 2018. An empirical analysis of journal policy effectiveness for
491        computational reproducibility. *Proceedings of the National Academy of Sciences of the United*
492        *States of America* 115:2584–2589. DOI: 10.1073/pnas.1708290115.

493    Taschuk M, Wilson G. 2017. Ten simple rules for making research software more robust. *PLoS*
494        *computational biology* 13:e1005412. DOI: 10.1371/journal.pcbi.1005412.

495    The data deluge. 2012. *Nature Cell Biology* 14:775–775. DOI: 10.1038/ncb2558.

496    Thiébaut A, Delaveau T, Benchouaia M, Boeri J, Garcia M, Lelandais G, Devaux F. 2017. The
497        CCAAT-Binding Complex Controls Respiratory Gene Expression and Iron Homeostasis in
498        Candida Glabrata. *Scientific Reports* 7. DOI: 10.1038/s41598-017-03750-5.

499    Ullah A, Lopes MI, Brul S, Smits GJ. 2013. Intracellular pH homeostasis in Candida glabrata in
500        infection-associated conditions. *Microbiology (Reading, England)* 159:803–13. DOI:
501        10.1099/mic.0.063610-0.

502    Vasilevsky NA, Minnier J, Haendel MA, Champieux RE. 2017. Reproducible and reusable research:
503        are journal data sharing policies meeting the mark? *PeerJ* 5:e3208. DOI: 10.7717/peerj.3208.

504

16

Legend:
- Required information *before* new Pixel Set submission (blue)
- Information detailed in the Excel file *during* new Pixel Set submission (yellow/gold)
- Information added through the Web Interface (purple)

**TagTreeModel** ‹BaseTagTreeModel,TagModel›

| parent | ForeignKey (None) |
|---|---|
| count | IntegerField |
| label | IntegerField |
| level | IntegerField |
| name | CharField |
| path | TextField |
| protected | BooleanField |
| slug | SlugField |

parent (children)

**Pixel**

| id | UUIDField |
|---|---|
| omics_unit | ForeignKey (id) |
| pixel_set | ForeignKey (id) |
| quality_score | FloatField |
| value | FloatField |

omics_unit (pixel)   pixel_set (pixel)

**Omics Unit**

**OmicsUnit**

| id | UUIDField |
|---|---|
| reference | ForeignKey (id) |
| strain | ForeignKey (id) |
| type | ForeignKey (id) |
| status | PositiveSmallIntegerField |

**Pixel Set**

**PixelSet**

| id | UUIDField |
|---|---|
| analysis | ForeignKey (id) |
| cached_omics_areas | ArrayField |
| cached_omics_unit_types | ArrayField |
| cached_species | ArrayField |
| description | TextField |
| pixels_file | FileField |

**SubmissionProcess**

| process_ptr | OneToOneField |
|---|---|
| analysis | FileField |
| archive | FileField |
| downloaded | BooleanField |
| imported | BooleanField |
| label | CharField |
| meta | JSONField |
| tags | JSONField |
| template_checksum | CharField |
| template_version | CharField |
| uploaded | BooleanField |
| validated | BooleanField |

type (omics_unit)   strain (omics_unit)

analysis (pixelset)   analysis (submissions)

**Omics Unit Type**

**OmicsUnitType**

| id | UUIDField |
|---|---|
| description | TextField |
| name | CharField |

**Strain**

**Strain**

| id | UUIDField |
|---|---|
| reference | ForeignKey (id) |
| species | ForeignKey (id) |
| description | TextField |
| name | CharField |

**Analysis**

**Analysis** ‹TaggedModel›

| id | UUIDField |
|---|---|
| pixeler | ForeignKey (id) |
| completed_at | DateTimeField |
| created_at | DateTimeField |
| notebook | FileField |
| saved_at | DateTimeField |
| secondary_data | FileField |

species (strain)

reference (omicsunit)

experiments (analysis)

pixeler (analysis)

**Pixeler**

**Pixeler** ‹AbstractUser›

| id | UUIDField |
|---|---|
| date_joined | DateTimeField |
| email | EmailField |
| first_name | CharField |
| is_active | BooleanField |
| is_staff | BooleanField |
| is_superuser | BooleanField |
| last_login | DateTimeField |
| last_name | CharField |
| password | CharField |
| username | CharField |

**Species**

**Species**

| id | UUIDField |
|---|---|
| reference | ForeignKey (id) |
| repository | ForeignKey (id) |
| description | TextField |
| name | CharField |

**Experiment**

**Experiment** ‹TaggedModel›

| id | UUIDField |
|---|---|
| omics_area | ForeignKey (id) |
| completed_at | DateField |
| created_at | DateTimeField |
| description | TextField |
| released_at | DateField |
| saved_at | DateTimeField |

reference (species)

reference (strain)

experiments (analysis)

omics_area (experiment)

tags (analysis)

tags (experiment)

groups (user)   user_permissions (user)

entries (experiment)

**Group**

**Permission**

**Entry**

**Entry**

| id | UUIDField |
|---|---|
| repository | ForeignKey (id) |
| description | TextField |
| identifier | TextField |
| url | URLField |

**Omics Area**

**OmicsArea** ‹MPTTModel›

| id | UUIDField |
|---|---|
| parent | ForeignKey (id) |
| description | TextField |
| level | PositiveIntegerField |
| lft | PositiveIntegerField |
| name | CharField |
| rght | PositiveIntegerField |
| tree_id | PositiveIntegerField |

**Tag**

**Tag** ‹TagTreeModel›

| id | AutoField |
|---|---|
| parent | ForeignKey (id) |
| count | IntegerField |
| label | IntegerField |
| level | IntegerField |
| name | CharField |
| path | TextField |
| protected | BooleanField |
| slug | SlugField |

parent (children)

repository (species)

repository (entry)

**Repository**

**Repository**

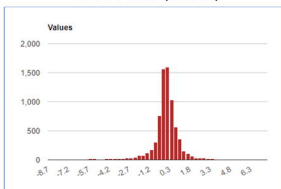| id | UUIDField |
|---|---|
| name | CharField |
| url | URLField |

**1 – Data exploration in the Pixel Web App**
Keywords: « Candida glabrata » and « alkaline pH »

✓ Pixel Set A
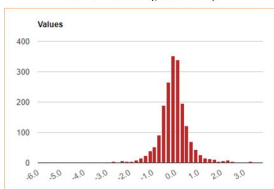Tags: Candida glabrata, WT, **transcriptomics**, alkaline pH, logFC

✓ Pixel Set B
Tags: Candida glabrata, WT, **proteomics**, alkaline pH, logFC

5253 Pixels (mRNA)

1879 Pixels (proteins)
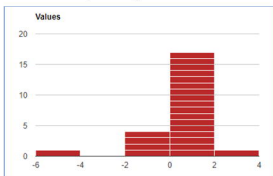
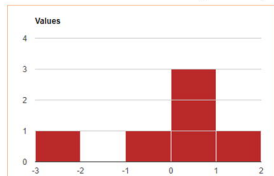**2 – Omics Unit filtering in the selected Pixel Sets**
Keywords: « pathogenesis »

| Omics Unit | Description |
|---|---|
| CAGL0C05467g | Ortholog(s) have role in biofilm formation, filamentous growth, pathogenesis and cytoplasm, nucleus localization |
| CAGL0D02156g | Ortholog(s) have glucosamine 6-phosphate N-acetyltransferase activity, role in UDP-N-acetylglucosamine biosynthetic process, pathogenesis and cytoplasm, nucleus localization |
| CAGL0F04807g | Ortholog(s) have role in pathogenesis and cell surface, hyphal cell wall, integral component of mitochondrial outer membrane, plasma membrane localization |

17 Pixels (mRNA)

6 Pixels (proteins)

**3 – Multi-pixel sets export for a new data analysis cycle**