

# GigaScience

## The on-premise data sharing infrastructure e!DAL: Foster FAIR data for faster data acquisition

--Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-20-00181	
<b>Full Title:</b>	The on-premise data sharing infrastructure e!DAL: Foster FAIR data for faster data acquisition	
<b>Article Type:</b>	Technical Note	
<b>Funding Information:</b>	Bundesministerium für Bildung und Forschung (031A053)	Dr. Daniel Arend
	Bundesministerium für Bildung und Forschung (031B0770A)	Mr. Patrick König
	Bundesministerium für Bildung und Forschung (031A536A)	Dr. Daniel Arend
<b>Abstract:</b>	<p>The FAIR data principle as a commitment to support long-term research data management is widely accepted in the scientific community. Although the ELIXIR Core Data Resources and other established infrastructures provide comprehensive and long-term stable services and platforms for FAIR data management, a notifiable number of research data is still hidden or under risk of getting lost. Currently, high-throughput plant genomics and phenomics technologies are producing research data in abundance, the storage of which is not covered by established core databases. Due to its high volume and strong heterogeneity, resulting in missing infrastructures these data sets often remain unpublished. To master these challenges the ELIXIR Germany/de.NBI service Plant Genomic and Phenomics Research Data Repository (PGP) implements a “bring the infrastructure to the data” approach which allows to keep the research data in place and wrap them by a FAIR-aware software infrastructure. This article presents new features of the PGP infrastructure as a best practice on how to easily set up FAIR-compliant and intuitive research data services. Furthermore, the integration of the ELIXIR AAI and data discovery services are introduced as means to lower technical barriers and to increase the visibility of research data.</p>	
<b>Corresponding Author:</b>	Daniel Arend, Ph.D. Leibniz-Institut für Pflanzengenetik und Kulturpflanzenforschung Seeland OT Gatersleben, SACHSEN-ANHALT GERMANY	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Leibniz-Institut für Pflanzengenetik und Kulturpflanzenforschung	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Daniel Arend	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Daniel Arend	
	Patrick König	
	Astrid Junker	
	Uwe Scholz	
	Matthias Lange	
<b>Order of Authors Secondary Information:</b>		
<b>Additional Information:</b>		

Question	Response
<p>Are you submitting this manuscript to a special series or article collection?</p>	<p>No</p>
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	<p>Yes</p>
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p>	<p>Yes</p>

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)



Placeholder for  
journal logo  
gigascience-  
logo.pdf

*GigaScience*, 2018, 1–10

doi: [xx.xxxx/xxxx](#)

Manuscript in Preparation

Technical Note

## TECHNICAL NOTE

# The on-premise data sharing infrastructure e!DAL: Foster FAIR data for faster data acquisition

Daniel Arend<sup>1\*</sup>, Patrick König<sup>1</sup>, Astrid Junker<sup>1</sup>, Uwe Scholz<sup>1</sup> and Matthias Lange<sup>1</sup>

<sup>1</sup>Leibniz Institute of Plant Genetics and Crop Plant Research (IPK) Gatersleben

\*arendd@ipk-gatersleben.de

## Abstract

The FAIR data principle as a commitment to support long-term research data management is widely accepted in the scientific community. Although the ELIXIR Core Data Resources and other established infrastructures provide comprehensive and long-term stable services and platforms for FAIR data management, a notifiable number of research data is still hidden or under risk of getting lost. Currently, high-throughput plant genomics and phenomics technologies are producing research data in abundance, the storage of which is not covered by established core databases. Due to its high volume and strong heterogeneity, resulting in missing infrastructures these data sets often remain unpublished. To master these challenges the ELIXIR Germany/de.NBI service *Plant Genomic and Phenomics Research Data Repository (PGP)* implements a “bring the infrastructure to the data” approach which allows to keep the research data in place and wrap them by a FAIR-aware software infrastructure. This article presents new features of the PGP infrastructure as a best practice on how to easily set up FAIR-compliant and intuitive research data services. Furthermore, the integration of the ELIXIR AAI and data discovery services are introduced as means to lower technical barriers and to increase the visibility of research data.

**Key words:** research data management; FAIR principles; digital object identifier; plant genomics and phenomics

## Introduction

The FAIR (Findability, Accessibility, Interoperability, and Reusability) principles, drafted by the FORCE11 workgroup in 2015 [1] and published in 2016 by Wilkinson et al. [2] are widely accepted and are increasingly adopted in the management policies of research data. The scientific community is showing a rising awareness of the scientific value of reusable research data. This has already resulted in FAIR principles being formally accepted in several data management guidelines, e.g. in the Horizon2020 program [3] of the European Commission, and integrated into the research funding policy [4, 5]. Its technical implementation is supported by data repositories, which store and share research data in a FAIR manner. Those can be classified into (i) general purpose data repositories, e.g. figshare [6], Zenodo [7] and Dryad [8], (ii) core data deposition databases, i.e. ELIXIR deposition databases for life science data [9], NCBI database resources [10] and (iii) institutional

databases and repositories. All have in common that the research data has to be transferred by its owner from the place of data generation to these repositories. This involves considerable effort for data compilation, cleansing, homogenisation, metadata enrichment, formatting and upload. As a result, the published datasets are condensed and generally limited to insufficiently documented supplement material for publications in scientific journals. In the case that data is prepared for a submission to database systems, e.g. the EBI and NCBI core data resources, Bioinformaticians are charged and trained. Examples are the preparation of data for submission to the EBI ENA archive [11, 12], the European Variation Archive (EVA) [13] or the preparation of ISA-TAB compatible data submission for plant phenotyping data [14, 15]. Alternatively, project-related or institutional data repositories could be set up. This in turn requires skilled technicians and computer scientists as well as long-term access to appropriate network and storage infrastructure. Such repositories frequently have a short life time,

Compiled on: July 6, 2020.

Draft manuscript prepared by the author.

whether due to staff fluctuation, long-term maintenance costs and resource consumption. Another reason may be that the repository's niche is too specific to attract substantial data volume, which in turn strongly depends on policies and cost-benefit considerations.

Thus, there is a need for an additional class of repositories that support the data sharing for this class of research data by moving the infrastructure to the data. The concept is to apply an on-premise, infrastructure-to-the-data (I2D) principle. The basic idea of the I2D approach is shown in Figure 1. In contrast to the conventional data publication pipelines to journal accepted databases, which usually involves a time-consuming data upload to an external platform and possibly additional costs depending on the required storage space, the underlying e!DAL software [16] encapsulates an existing storage infrastructure by a data publication layer. This layer is a broker to the DataCite [17] data publication service agent and provides an API and a tooling infrastructure for data submission, DOI delivery, reporting and data quality reviewing. This finally enables the assignment of DOIs with a minimal set of metadata to in-house stored data and its approved FAIR referencing by journals or data lookup services.

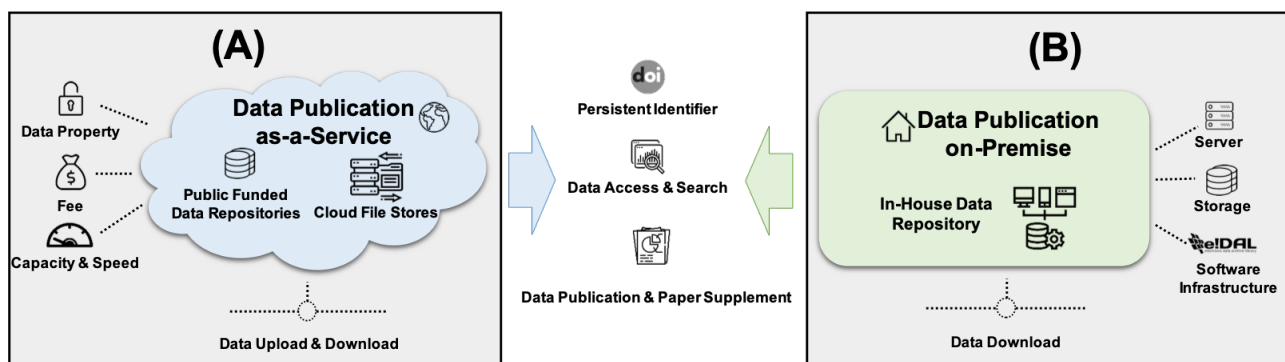
As proof of concept, the Plant Genomics and Phenomics Data Repository (PGP) was implemented [18] to publish digital plant genetic resources (PGR) [19] according to the FAIR principles. Approximately seven million crop accessions are conserved in genbank collections worldwide. The valorisation of PGRs through genotyping and phenotyping is of special focus in the public and private sectors [20, 21]. The data management of digital PGRs is identified as one of the most important challenges for a long-term strategy to enhance the productivity, sustainability and resilience of crop varieties and agricultural systems. In contrast to successful studies on genomics-assisted genbank management and the utilization of germplasm collections [20], the special focus of the PGP repository is the publication of buckets of research data that do not fit into classical databases due to their volume, objective, structure or incomplete analysis. Examples are primary data from imaging, field phenotyping, SNP matrices, 3D-plant models, metabolite screenings and environmental sensor data. The experience gained during the four-year operation of the repository has led to a growing acceptance of this approach for the publication of digital PGRs collected in the context of the German Federal ex situ Genebank of Agricultural and Horticultural Crop Species [22]. This experience and the adoption as a service in the European life-sciences Infrastructure for biological Information ELIXIR [23] resulted in novel features, which were implemented with the aim of further improving its acceptance and enabling increased sharing of digital PGRs. After an

update to the state of the art, the new features of the e!DAL data sharing software and its application for the publication of digital PGRs will be explained.

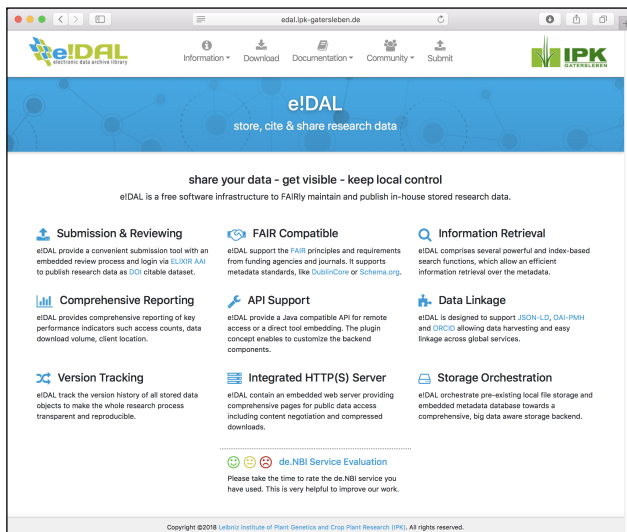
## Related Work

Just like we have many different data types from several domains, there are also a variety of domain-specific archives and information systems. Most of them evolved over many years and they are widely accepted by the research community [24], *ENA* for genomic data [25], *UniProt* for protein data [26], *PRIDE* for proteome data [27], *BioModels* for system biology data [28] and many more. Usually also research journals and other publishers request data sharing using these established domain-specific databases. In order to not get lost in the diversity of archives, there are several registries like *re3data.org* or *FAIR-sharing.org*, as well as consortia like *GFBio* which collect and categorize repositories to help researchers finding the suitable storage for their data.

Infrastructure programs like the European Open Science Cloud (EOSC) [4], and the European life-sciences Infrastructure for biological Information (ELIXIR) [29] coordinate maintenance and interoperability of research data repositories as federated services by member organisations and hosting institutions. Furthermore, the ELIXIR organisation has the aim to establish a stable and sustainable infrastructure for biological information. In doing so, they define important core resources and tools as a support of the research community [9] like *BRENDA* [30] or *SILVA* [31]. Most of these systems accept only very specific datasets and require specialised metadata based on schemes that have been improved by the community over years. Unfortunately, there are several, mostly relatively new data types, e.g. plant phenotypic data, which currently do not fit into existing databases, mainly because of their strong heterogeneity and high volume. Public data sharing services like *figshare* or *DRYAD* provide an alternative solution for publishing these datasets. They are easy to use and have a comprehensive functionality like supporting version controlling and the assigning of persistent identifiers. One important deficiency of such services is the limited free space, which is usually enough for sharing some reduced graphics or aggregated tables, but not for storing large datasets. Furthermore, the establishment and configuration of an own in-house infrastructure based on existing software packages like *CKAN* or *Dataverse* could overcome this shortcoming, but they require quite a lot of technical prerequisites and know-how.



**Figure 1.** Data-Publication-as-a-Service vs. Data-Publication-on-Premise: Both services feature FAIR data publication, but differ in costs. The as-a-service model (A) costs fee, delegation of data property control and face capacity limits in storage and data upload. The on-premise model (B) keep data in-house, but requires the availability of server and storage hardware as well as the installation of the e!DAL software.



**Figure 2.** Screenshot of the e!DAL project website. On the project website a very detailed description on how to use the e!DAL as well as comprehensive examples for users as well as for developers. Furthermore some short video tutorials and webinar recordings are available to lower the initial barriers to dig into the infrastructure

## Infrastructure

To lower the technical barriers and minimize the effort for scientists to archive and share their research data, we developed the generic **e!DAL software infrastructure** [16]. The usual “Data Publication as-a-Service” procedure includes the transfer of selected datasets to external databases and storage infrastructures after data generation and analysis. In this way research data can be referenced in a future research publication, as show on the left side of Figure 1 (A). In contrast the e!DAL infrastructure provide a “Data Publication on-Premise” approach which enables the publication of locally stored high voluminous research data through the assignment of widely accepted and long-term stable Digital Object Identifiers (DOIs). This is illustrated on the right side of Figure 1 (B). Using DOIs for referencing provides multiple advantages for sharing and accessing research data. Beside adding them as supplements to a research article they can also be the basis for a comprehensive data paper [32]. Furthermore, the well connected infrastructure of the DataCite consortia strongly increase the visibility of the research data assigned with a DOI. It is automatically linked with the ORCID account of the authors, can be found via the DataCite Search and other common search engines or can be harvested via OAI-PMH interface.

Based on the e!DAL infrastructure, the **Leibniz Institute of Plant Genetics and Crop Plant Research (IPK) Gatersleben** and the **German Plant Phenotyping Network (DPPN)** jointly have initiated the **Plant Genomics and Phenomics Research Data Repository (PGP)** [18] as a powerful infrastructure for the publication of comprehensive plant genomics and phenomics research data. The repository covers in particular cross-domain datasets, which are not being published in public repositories for reasons of data volume or data domain, such as phenotyping images, genotyping data, visualizations of morphological models, data from mass spectrometry as well as software and related documents. Doing so, PGP currently provides 200 data records, which can be referenced via DOIs and are annotated with technical metadata. These records comprise more than 1.4 million files with an overall volume of over 2.6 terabytes (see Figure 3). To ensure data discoverability, PGP provides landing pages with JSON-LD formatted metadata and is therefore discoverable through data web crawler services, which follow the

schema.org recommendations, such as Google, Microsoft, Yandex etc. Furthermore, e!DAL implements the OAI-PMH Protocol for Metadata Harvesting from Open Archives Initiative. To support scientists to disseminate their research data the PGP infrastructure is accepted as institutional repository for the **Scientific Data** (Nature Publishing Group) and the **GigaScience** (Oxford Academic) Journals and is registered in [re3data.org](https://re3data.org), [FAIRsharing.org](https://fairsharing.org), [OpenAIRE](https://openaire.org) and [DataCite](https://datacite.org).

The benefits of this wide support of data discovery enabling technologies and data publication in general is proven by the steadily increasing number of dataset accesses. By June 2020, PGP delivered 300 terabytes of data and the provided datasets have been accessed by 100,000 unique clients.

## Improvements

The following section will sum up the main enhancements and updates of the e!DAL infrastructure which comprise new general features, comprehensive changes of several frontend components and important performance improvements. Furthermore, an extensive update due to the latest changes in the Java programming languages and an improved build and deployment process is described.

## Performance

After releasing the first productive version of the PGP repository in 2015, we received many diverse data submissions from several research domains and with very heterogeneous data files. Since then we recognized that the e!DAL infrastructure software scales very well and is able to handle millions of data files, which confirms previous calculations and performance tests [16]. But furthermore it became apparent that sometimes the performance decreases, e.g. for uploading comprehensive datasets with several hundred thousands of small files. Since this is a very common case, e.g. for plant phenotyping datasets, an improvement of the implementation of the e!DAL infrastructure was necessary.

One important feature of e!DAL is the automatic calculation of several essential technical metadata, like the MIME-Type, the data volume or the checksum of every file when storing new datasets. This is convenient, because the user does not need to provide these information on his/her own, but of course these computations are a resource- and runtime-intensive process. Therefore, the functionality to determine the previously mentioned technical metadata and the procedure to transfer the actual binary data have been improved towards a parallel processing of multiple files. This results in a better performance especially on today’s multi-core systems. Furthermore, we optimized several settings for the streaming buffer size and the remote transfer to improve the memory usage and the upload performance for the case of numerous small files. Additionally the checksum calculation was updated to use the more collision resistant SHA-256 algorithm, instead of the older and insecure MD5 function.

## New Features

The previous version of the e!DAL infrastructure already fulfilled several recommendations of the FAIR data principles, such as the support of standardized metadata based on the DublinCore schema or the provision of persistent DOIs for accessing and referencing of research datasets. The e!DAL infrastructure has been further updated to optimize the usability and the general user experience. Additional features were implemented to increase the visibility of published data and the

acceptance of the infrastructure, which in the end also led to be even more FAIR compliant. Thereby the roadmap for scholarly data repositories [33] was taken into account. Subsequently, the most important extensions are described.

### ORCID

To efficiently find and access specific research data files across millions of datasets, persistent identifiers like DOIs or URNs are very helpful and well established. But nevertheless, the research community is also quite large, and sometimes it is very difficult to distinguish data authors because of similar names or to identify the same researcher after he changes his affiliation. With the *Open Researcher and Contributor ID (ORCID)* there is an easy and persistent solution to uniquely identify authors and to solve issues with name ambiguity [34]. An important advantage is the interdisciplinarity, because ORCID is used across nearly all research domains and organizations, e.g. at the mid of 2019 in Germany there were already 150.000 ORCID's registered [35]. By linking authors with publications, affiliations or funding agencies, it helps to find relationships between researchers and their work and the corresponding research data. Since the e!DAL infrastructure is generic and suitable for different kinds of research data, the ORCID system gives us an ideal solution to identify authors and improve the collected metadata for published datasets. Furthermore, the authors and their research data will get a better visibility, due to the connection between the ORCID infrastructure and infrastructure of the DataCite consortium, which is handling the DOIs.

To add the possibility for assigning an ORCID to every data creator or contributor in the e!DAL infrastructure, the PERSON data type in the metadata schema was extended. e!DAL uses the REST API of the ORCID registry to provide the possibility to search for the ORCID of a given name. In addition, it can be validated if an entered ORCID belongs to the corresponding name to prevent an accidentally linking with a wrong ORCID. All these API functions were integrated into the graphical user interface of the data submission tool for the PGP repository. Furthermore, the content pages of published and DOI linked datasets were improved to provide direct links to the ORCID profiles of the associated authors and contributors of the data.

### JSON-LD & DC meta tags

Another method of making research data interoperable as well as machine-readable is to embed the describing metadata using *JavaScript Object Notation for Linked Data (JSON-LD)* format. This approach provides comprehensive possibilities to harvest

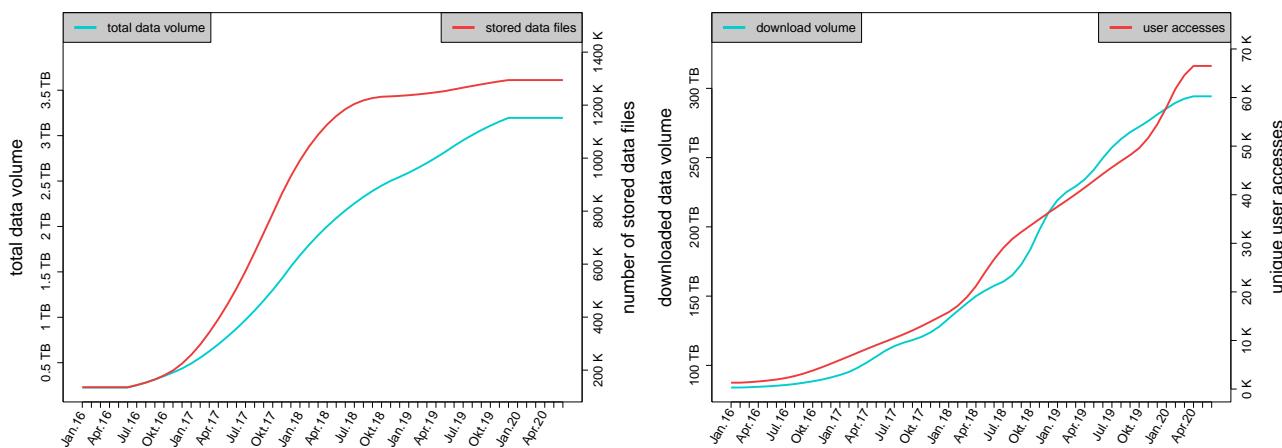
and reuse research data. JSON-LD is a data serialization and exchange method and was developed to be easily embeddable into various systems for providing interoperable web services [36]. The dynamic HTML templates for the content pages of the embedded webserver of e!DAL, which provides the URLs for resolving the assigned DOIs, have been extended accordingly.

```
<script type="application/ld+json">
{
  "@context": "http://schema.org",
  "@type": "Dataset",
  "@id": "https://doi.org/10.5447/IPK/2016/7",
  "name": "Raw images files from quantitative monitoring of...",
  "publisher": {
    "@type": "Organization",
    "name": "IPK Gatersleben"
  },
  "description": "This dataset contains 30426 raw image files...",
  "keywords": "high throughput plant phenotyping, growth protocol...",
  "inLanguage": "en",
  "author": {
    "@type": "Person",
    "givenName": "Astrid",
    "familyName": "Junker",
    "address": "IPK Gatersleben"
  },
  "contributor": {
    "@type": "Person",
    "givenName": "Thomas",
    "familyName": "Altmann",
    "address": "IPK Gatersleben"
  }
}
```

**Listing 1.** Reduced example of the JSON-LD data from the content page of a DOI assigned with e!DAL, which is stored in the PGP repository

Listing 1 show an example for the JSON-LD description of a dataset in the PGP repository. The attributes are based on the [schema.org](http://schema.org) ontology, which is a well-established and community-driven vocabulary used to structure digital data on websites. It is used and harvested by several common search engines [37] and provides an interoperability between dataset from separated resources and platforms.

Another alternative to JSON-LD are so-called HTML meta tags. They are embedded in the <head> section of an HTML-document and also allow to harvest the metadata and describe connections between datasets from different infrastructures. As the metadata schema of the e!DAL infrastructure is already inspired by the DublinCore metadata schema [38] the embedded HTML templates for the content pages of published datasets were extended to provide the technical metadata of every object also as HTML meta tags (see Listing 2).



**Figure 3.** Access numbers and stock overview of the eDAL-PGP repository since 2016. The left diagram shows the development of the data volume and number of stored files, which were provided by the repository. The diagram on the right side the constantly increasing number of accesses and downloads of published datasets is shown.

```

<meta name="DC.Title" content="Screening_of_wild_potato_genetic...">
<meta name="DC.Identifier" content="https://doi.org/10.5447/IPK/2019/1">
<meta name="DC.Publisher" content="e!DAL_-_Plant_Genomics_and_Phenomics...">
<meta name="DC.Language" content="en">
<meta name="DC.Description" content="This_data_set_contains_results_of...">
<meta name="DC.Rights" content="CC_BY-NC-SA_4.0">
<meta name="DC.Creator" content="Bachmann-Pfabe,_Silvia...">
<meta name="DC.Contributor" content="Dehmer,_Klaus...">
<meta name="DC.Subject" content="Phytophthora_infestans">
<meta name="DC.Subject" content="germplasm_collection">

```

**Listing 2.** Reduced example of the DublinCore Meta-Tags from the content page of a DOI assigned with e!DAL

### Content Negotiation

Persistent DOIs provide a solution for long-term stable resolvability and referencing of all published datasets. In addition, for several reasons such as citing the datasets or harvesting the metadata, it is necessary to provide content negotiation to serve resources in different formats. Therefore the possibility to get different representations of the public datasets stored in an e!DAL infrastructure was implemented and can be used by several export functions, which were added on the corresponding content pages as shown in Figure 4. They provide the option to get textual representations, citation formats like BibTeX or RIS and linked data formats like schema.org/JSON-LD and RDF for every dataset. Due to the fact that the [DataCite service](#) already provides a content negotiation feature, it was not necessary to implement a separate function for the embedded webserver of e!DAL. Instead, the HTTP handler uses the provided function for the different formats via a REST call and redirects the responses to the e!DAL infrastructure.

The screenshot shows a web browser displaying a dataset page from the e!DAL PGP Repository. The page header includes the e!DAL logo and the IPK Gatersleben logo. The main content area features a citation for 'A. Junker (2016-02-17): Image analysis results from quantitative monitoring of 484 Arabidopsis thaliana plants using high-throughput plant phenotyping. DOI:10.5447/IPK/2016/9'. Below the citation, a red rectangle highlights a row of download buttons: BibTeX, RIS, RDF/XML, RDF/Turtle, Text, and JSON-LD. The abstract text describes the dataset's content, and the license is listed as CC BY 4.0. The DOI and content size (90 Directories, 29725 Files, 18.5 GB) are also visible.

**Figure 4.** Screenshot of an example datasets in the PGP repository. The red colored rectangle mark the e!DAL embedded functions for the content negotiation.

### Elixir AAI

The e!DAL infrastructure provides a flexible and embedded security concept based on the Java Authentication and Authorization Service (JAAS). To provide the research data management and publication capabilities to a wide range of users from universities, research institutes or further organisations, a new login module using the ELIXIR Authentication and Authorization Infrastructure (AAI) [39] was implemented. It was designed to provide a single sign-on service for authenticating researchers to services, which are a part of the ELIXIR portfolio. Doing so, it combines the huge amount of existing organisational identity providers from institutes that are associated with ELIXIR under one roof.

The screenshot shows a web browser displaying the ELIXIR AAI login dialog. The page title is 'Select your identity provider'. There are three main buttons: 'Sign in with Google' (red), 'Sign in with LinkedIn' (blue), and 'Sign in with ORCID' (green). Below these, there is a section for 'your institutional account' with a search box containing 'Jülich' and a dropdown menu showing 'Forschungszentrum Jülich GmbH'. At the bottom, there is a link 'Can't find your institution? Select it in extended list and help us' and a button 'Add your institution'.

**Figure 5.** Integration of the ELIXIR AAI login dialog into the e!DAL infrastructure. The ELIXIR AAI provides beside a large collection of institutional identity providers, which can be easily found via the small search box, also several third party login opportunities based on Google, ORCID or LinkedIn. This provides the users a very flexible authentication infrastructure.

The new e!DAL login module follows the OAuth protocol [40] to authenticate users over the ELIXIR AAI and automatically receive their email address, which is necessary for the communication between the data submitting researcher and the reviewers of the embedded review process. As the first use case, the new ELIXIR AAI based login was integrated into the PGP repository to open the infrastructure and the data submission process for offering the service to a wide range of researchers without the need of creating a separate account. The ELIXIR AAI allows researchers to use their existing organisational accounts (see Figure 5), which lowers the barrier to use the infrastructure and to reach a larger group of data providers.

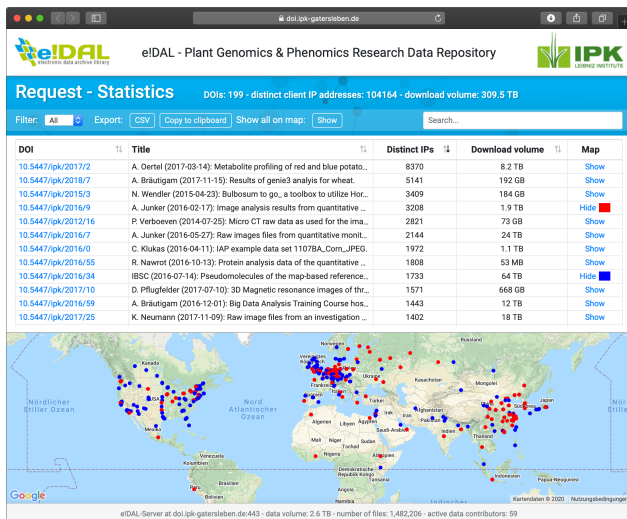
Furthermore, with the opportunity to use the ELIXIR AAI, the already low effort, which is necessary to establish further e!DAL installations, was reduced. Therefore at the end of 2018 a further e!DAL based repository at the [Jülich Plant Phenotyping Center \(JPPC\)](#) was established using the ELIXIR AAI login provider.

### Amended frontend

The [Apache Velocity](#) template engine is used to render all HTML-based content of the e!DAL embedded webserver like the landing pages of published datasets and e-mail messages. This prevents the infrastructure from storing a massive amount of very similar websites and text drafts, which saves storage and provides a high performance for delivering content via the HTTP handler. All websites are provided dynamically on demand and created from only a few reusable templates.

For the latest e!DAL version all content pages and the underlying templates were fully redesigned to provide a pleasing visual look and functional user experience. By using frontend frameworks and libraries like [Bootstrap](#) and [jQuery](#) it is ensured that the user interface is responsive and working on both modern desktop browsers as well as on mobile devices. Figure 6 shows the new layout as an example screenshot of the embedded access statistic page of the PGP repository. Together with the new design for the frontend components of e!DAL, also the project website was renewed to provide comprehensive infor-





**Figure 6.** Example screenshot of the e!DAL embedded report page for the e!DAL-PGP repository showing the new layout and user interface. Several simple filter and search possibilities help to look for specific datasets. An integrated world map gives the user an idea on the dissemination of the data.

mation for the user and for developers in a more concise manner.

## Deployment and Usability

Since the last major release of the e!DAL infrastructure software a lot of optimizations and several new functionalities, which were described in the previous sections, have been implemented. Together with these improvements, changes in the general build and release process and in the usability have also been integrated. The most relevant of them are explained subsequently.

### Gradle Multi-Build Project

After using the [Maven](#) build system Build System for several years for developing and releasing the e!DAL software components a change to the [Gradle](#) build tool was performed. Due to the constant increasing size of the project and the source code, because of new functionalities, several extensions and additional unit tests to guarantee a high software quality the build process using Maven takes quite a long time. This makes the regular release of stable versions very time-intensive. Furthermore, the build configuration became more complex and difficult to maintain. Gradle is strongly focused on a fast and specific build cycle. It supports multi-core systems to a high degree and allows e.g. the execution of several test suites in parallel. With the change to the build infrastructure, we also decided to redesign the entire project build hierarchy and created a multi-build project for the e!DAL infrastructure. It contains the main API components including the reference implementation as well as the components for the server-client architecture, which is directly based on this core implementation. This approach massively accelerates the build time, simplifies the maintenance and allows a more frequent deployment of new versions. The project is now available in a new [BitBucket repository](#).

Nevertheless, the API is still released as an artifact in the central Maven Repository and can be integrated into other software projects using Maven or Gradle, as shown in Listing 3.

```
----- Maven - 'pom.xml' -----
<project>
<dependencies>
<dependency>
<groupId>de.ipk-gatersleben</groupId>
```

```
<artifactId>eDAL-MetaDataAPI | eDAL-MetaDataAPI-Server | eDAL-
MetaDataAPI-Client</artifactId>
<version>3.0.2</version>
</dependency>
<dependencies>
</project>
```

```
----- Gradle - 'build.gradle' -----
repositories {
mavenCentral()
}
dependencies {
compile 'de.ipk-gatersleben:eDAL-MetaDataAPI:3.0.2'
compile 'de.ipk-gatersleben:eDAL-MetaDataAPI-Server:3.0.2'
compile 'de.ipk-gatersleben:eDAL-MetaDataAPI-Client:3.0.2'
}
```

**Listing 3.** Integration of the e!DAL components into the configuration of an Apache Maven or Gradle based project

### OS specific executables

Due to the complete new development and release cycle by Oracle, the [Java](#) programming environment, which is the basis for the e!DAL infrastructure, changes a lot in recent years. In addition, the comprehensive redesign and reconstruction of the language itself, like the introduction of the new module concept or the removal of popular and formerly native APIs and frameworks like [JavaFX](#) or the [Java Network Launching Protocol \(JNLP\)](#) were some very substantial changes. This strongly influences the e!DAL implementation, because they were also a significant part of the previous version. Unfortunately this impeded at some points the further development of e!DAL infrastructure, because a lot of the used frameworks and libraries needed several months to update their code to be compatible with the latest Java versions. With the new version 3.0.0 the e!DAL infrastructure is fully based on the [Java Runtime Environment \(JRE\) 12](#). Therefore some comprehensive changes were necessary. In order to run e!DAL with the different existing runtimes, e.g. the official runtime from Oracle, but also the alternative and widely used [OpenJDK](#), it was necessary to integrate the [JavaFX](#) library directly into the implementation. This increases the actual size of the API package, but it makes the infrastructure much more compatible and even more independent from the system preconditions than before.

The removal of the support for the popular and well known [JNLP](#) was also a high challenge, because the [Java webstart](#) tool was used to give the user an intuitive and platform-independent way to run the graphical data submission tool. Nevertheless this solution also provides some shortcomings like the need to provide an installed and compatible Java runtime. With the recently developed [jpackage](#) Java provides a powerful tool to pack self-contained applications along with a suitable JRE. We used [jpackage](#) to create a full image of the e!DAL data submission tool together with a reduced JRE, which contains only the necessary java modules and provides separate executables for the most common operation system (Windows, Unix, MacOS). This provides a very convenient usability for data submitter and makes the infrastructure again more compatible and independent from the given system preconditions of the users.

### Web-based submission application

In parallel to the update process due to the previously mentioned changes in the [Java Runtime Environment](#) and the development of the build process to create the self-executable applications for the submission dialog, a new web-based application was implemented to provide an alternative opportunity to upload research data to an e!DAL based infrastructure. The goal was the deployment of a user-friendly web application with the similar functionality of the corresponding desktop tool, but without the need to download the application as an executable or additional plugins. The [Vaadin](#) framework for Rich Internet Applications (RIA) was used for the implementation. [Figure 7](#)

**Plant Genomics & Phenomics Research Data Repository**  
Data Submission (BETA)

**Title**  
Image analysis results from quantitative monitoring of 484 *Arabidopsis thaliana* plants

**Description**  
This dataset contains image analysis results based on 30426 raw image files from a high-throughput plant phenotyping experiment. Here, the growth and development of 484 *Arabidopsis thaliana* plants (C24 genotype) in respect to plant movement and soil cover has been monitored over 55 days. This dataset contains 29724 feature-annotated result images. The features are plant height and volume.

**Authors**

FirstName	LastName	ORCID	LegalName	Address	ZIP	Country	Type	Remove
Daniel	Arend	0000-0002-2455-5938	Leibniz Institute of Plant Genetics and Crop Plant Research		06466	Germany	Creator	
Astrid	Junker	0000-0002-4656-0308	Leibniz Institute of Plant Genetics and Crop Plant Research		06466	Germany	Creator	

**License**  
 CC0 1.0 Universal (Creative Commons Public Domain Dedication)  
 CC BY 4.0 (Creative Commons Attribution)  
 CC BY-SA 4.0 (Creative Commons Attribution-ShareAlike)  
 CC BY-ND 4.0 (Creative Commons Attribution-NoDerivatives)  
 CC BY-NC 4.0 (Creative Commons Attribution-Non-Commercial)  
 CC BY-NC-SA 4.0 (Creative Commons Attribution-Non-Commercial-ShareAlike)  
 CC BY-NC-ND 4.0 (Creative Commons Attribution-Non-Commercial-NoDerivatives)

**Embargo Date**

**Next**

**Figure 7.** Example screenshot of the web-based submission application for the e!DAL-PGP repository. The form like user interface is very intuitive and similar to the desktop application. It helps the user to fill in the needed metadata by executing several checks to guarantee the completeness of the metadata.

shows a screenshot of the web application. By using several REST APIs, e.g. from the ORCID Registry or the ELIXIR AAI, a light-weight application could be created providing the same functionality as the full desktop client. Furthermore, users now have the possibility to submit research data also from mobile devices or other browser compatible devices. The only small shortcoming of the data submission via the web application is currently that not all browsers support the upload of comprehensive file folders. The latter is only possible if a recent version of Google Chrome or Mozilla Firefox is used. Other web browsers only allow the upload of single files.

## Results

In this article the basic overall ‘on-premise’ data management and publication concept of the e!DAL infrastructure as well as several new features and technical developments were presented. As a result, e!DAL matured to a comprehensive and FAIR-compliant infrastructure, while always keeping the focus on simple and flexible setup and integration into existing infrastructures and into the daily research process. With the described ‘bring the infrastructure to the data’ approach, it differs fundamentally from generic publication platforms like figshare or DRYAD, which can produce, depending on the needed storage, considerable financial costs and time costs for transferring the data. e!DAL allows the usage of available in-house storage capacities, without the need of complex requirements and technical infrastructures or comprehensive adaptations. All functionalities are already included and the provided reference implementation contains required components, such as a database or a webserver. This is a crucial advantage in comparison to other similar software infrastructures, like DataVerse or CKAN, and lowers the barrier to establish a publication infrastructure even for small-size research institutions with limited possibilities and know-how.

Thereby the FAIR compliance can be fulfilled by several e!DAL functions and components:

- **Findable:**

By providing embedded and machine-readable metadata based on standardized established formats, the e!DAL published datasets can be easily found using common search engines like Google or the DataCite Metadata Search. Due to

the widely established and used DOIs, the DataCite consortium is also involved in several projects and interacts with different systems like ORCID, CrossRef or Scholix. This further improves the findability of e!DAL datasets.

- **Accessible:**

The usage of DOIs as persistent identifiers guarantees a long-term stable availability of published datasets. The DataCite resolver for the DOIs allows simple access to the data and reference datasets, e.g. in a research article or as part of data publication. If the storage location of the underlying data is changing, the corresponding DOI remains stable and allows the uninterrupted access to the data by updating the resource path. Even if the data is no longer available the metadata is still accessible via the DataCite registry.

- **Interoperable:**

To provide interoperable datasets and to allow the aggregation of information about the relationship of datasets from different sources, the e!DAL infrastructure supplies embedded metadata on the content pages of every data object. They are stored using standardized formats and vocabularies like JSON/LD or rather schema.org.

- **Reusable:**

By collecting a standardized set of mainly technical metadata e!DAL guarantees a long-term readability and usability of all published datasets. The schema is inspired by the DublinCore metadata format and meets community established standards. Furthermore, a clear and easy license handling allows to assign a suitable license, which defines by whom and how the data can be used. They are available both on the content page of every data object as well as embedded in the HTML sources.

### e!DAL Usage

Established in 2016, the PGP repository is the first productive repository based on the e!DAL infrastructure and a part of the service portfolio of the GCBN unit (German Crop BioGreenformatics Network) [23] of de.NBI (German Network for Bioinformatics Infrastructure) [41], which is the head of ELIXIR Germany. After more than three years of productive usage, the PGP repository currently shares 200 comprehensive, plant-related research datasets containing mainly genomic and phenomic information, but also metabolic datasets or software components and pipelines. They have a total data volume of 2.6 terabytes consisting of more than 1,4 million files. Most of the datasets are part of a corresponding research paper and allow authors from IPK, but also from other institutes, to improve their manuscripts by enriching them with the underlying research data in a FAIR compliant way. Until June 2020 an overall download volume of around 300 terabytes and over 100,000 distinct user accesses were counted, which shows the high visibility of the provided data and the interest of the research community for this kind of data.

The integration of the ELIXIR AAI into the login mechanism of the PGP Repository is a prime example that shows how established platforms can benefit from the ELIXIR network. The provided services contribute to the increase of the visibility, to overcome the obstacles for the use of available infrastructures and to support FAIR compliant access to research data. The support of the ELIXIR single sign-on service enables collaborators to use the PGP repository as a service to publish their research data. Furthermore, the ELIXIR AAI login is fully integrated into the e!DAL infrastructure software, which allows to set up further FAIR in-house repository instances, following the presented I2D approach. Doing so in June 2018 a second repository based on the developed e!DAL infrastructure was established at the *Forschungszentrum Jülich*. Due to the auto-configuring installation it was possible to run the system and provide the submission and review workflow with only a

little effort in time. The integrated ELIXIR AAI login allows researchers from Jülich to use their existing institutional accounts. The complete infrastructure is hosted and maintained by the [Jülich Plant Phenotyping Center \(JPPC\)](#). The process of establishing further e!DAL-based repositories at the [Julius-Kühn Institute](#) and the [Helmholtz Centre München](#) are currently underway.

## Outlook

In this work, we showed the newly designed I2D concept for FAIR compliant data publication by using in-house storage infrastructures and new features of the e!DAL platform. After several years of operating a productive instance of this infrastructure as the basis for the PGP repository, we recorded high numbers of accesses and downloads. Although researchers have more and more possibilities to share their research data with the community, the incentive to do so is still not high enough for some researchers [42]. In contrast to the common peer-reviewed publication in journals, it is not so easy to measure the impact of research data itself, because the concept of data citation is still not a common practice [43], but it becomes more and more important and accepted [44]. It's not only a cultural problem, but also a technical challenge and therefore an issue of practicability [45]. One of the first metrics to count data citations was the commercial [Data Citation Index](#). But in the meanwhile some free and community initiated projects like [Make Data Count](#) have been developed. Furthermore, popular journals are starting to demand that authors put their research data as data citations in their common reference list [46]. This facilitates to measure its impact through a citation index and improves the visibility to readers which in turn increases the general acceptance of research data as valuable scientific assets. In future we will investigate several approaches for counting data citations and getting more credit for publishing research data. We plan to integrate a generic and open-source solution into the e!DAL infrastructure to show users comprehensive information how their data is reused and referenced.

The ORCID provide a widely accepted and used solution to unambiguously identify researchers. The integration within the e!DAL infrastructure is very intuitive and facilitates handling of multiple ORCIDs for comprehensive lists of authors. Besides the identification of persons, it can be also quite challenging to handle the diverse affiliations of research institutes, universities or companies with a focus on different scientific topics. Some authors have multiple affiliations, from time to time organizations may be renamed, the official addressee may change due to infrastructural developments or it may happen that an institute will be closed. The [Research Organization Registry \(ROR\)](#) provides an open and sustainable approach, which is led by the community and supported by popular organizations like DataCite or Dryad. The concept of the ROR identifiers is very similar to the ORCIDs and allows to uniquely identify all kinds of research organizations. Therefore, one of the next functional improvements for the e!DAL infrastructure will be the integration using the provided ROR API. This will cause some changes in the basic data structure, which however will result in a much easier and FAIRer way to handle author affiliations [47].

## Availability of source code and requirements

- Project name: e!DAL (electronic Data Archive Library)
- Current version: 3.0.2
- Project homepage: <https://edal.ipk-gatersleben.de>
- Source Code Repository: [https://bitbucket.org/ipk\\_bit\\_](https://bitbucket.org/ipk_bit_)

[team/electronicdataarchivelibrary](https://bitbucket.org/ipk_bit_team/electronicdataarchivelibrary)

- Operating system(s): platform independent
- Programming language: JVM based (Java 12+)
- JavaDoc: <https://edal.ipk-gatersleben.de/javadoc>
- Artifact repository: Maven Central (<https://mvnrepository.com/artifact/de.ipk-gatersleben>)
- License: GNU General Public License (GPL) Version 3 (<https://www.gnu.org/licenses/gpl-3.0.html>)

The e!DAL [project website](#) provide comprehensive information, which are relevant for users as well as developers. Beside a full Java documentation, further presentations, videos and several code and usage examples are presented. We changed the licence model of e!DAL to GNU General Public License (GPL) Version 3. This aim at maximising the spread of the e!DAL infrastructure in scientific community to foster FAIR principles of in-house stored data and to enable the incorporation of e!DAL into 3rd party software as well.

## Declarations

### List of abbreviations

AAI: Authentication and Authorization Infrastructure;  
 API: Application programming interface;  
 de.NBI: German Network for Bioinformatics Infrastructure;  
 DOI: Digital Object Identifier;  
 DPPN: German Plant Phenotyping Network;  
 FAIR: Findable, Accessible, Interoperable, Reusable;  
 GCBN: German Crop BioGreenformatics Network;  
 I2D: Infrastructure-to-the-Data (I2D);  
 JNLP: Java Network Launching Protocol;  
 JPPC: Jülich Plant Phenotyping Center;  
 JRE: Java Runtime Environment;  
 JSON-LD: JavaScript Object Notation for Linked Data;  
 ORCID: Open Researcher and Contributor ID;  
 PGP: Plant Genomics and Phenomics Research Data Repository;  
 PGR: Plant Genetic Resources;  
 RIA: Rich Internet Application;  
 RMI: Remote Method Invocation;

## Competing Interests

The authors declare no conflict of interest.

## Funding

This work was supported by the German Federal Ministry of Education and Research (BMBF) in frame of the the projects German-Plant-Phenotyping Network - DPPN (FKZ 031A053), Modernste Virtualitäts- und erweiterte Realitäts-Verfahren für den Zyklus von Samen zu Samen - AVATARS (FKZ 031B0770A) and German Network for Bioinformatics Infrastructure - de.NBI (FKZ 031A536A).

## Author's Contributions

Conceptualization: DA,ML  
 Software: DA,PK  
 Investigation: DA  
 Supervision: ML  
 Writing original draft: DA,ML  
 Writing review & editing: All  
 Funding acquisition: US,AJ,ML

## Acknowledgements

We want to thank Thomas Münch, Jens Bauernfeind and Heiko Mieke for the technical support. Furthermore we are very thankful to Benjamin Bruns, Robin Kluth and Michael Seidel for fruitful discussions, testing the infrastructure and formulating feature requests. Last but not least we thank all users and data producers for their valuable feedback.

## References

- Martone ME. FORCE11: Building the Future for Research Communications and e-Scholarship. *BioScience* 2015;65(7):635. <http://dx.doi.org/10.1093/biosci/biv095>.
- Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, et al. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific data* 2016;3. <http://dx.doi.org/10.1038/sdata.2016.18>.
- European Commission – Directorate-General for Research and Innovation. Guidelines on FAIR Data Management in Horizon 2020; 2016, [http://ec.europa.eu/research/participants/data/ref/h2020/grants\\_manual/hi/oa\\_pilot/h2020-hi-oa-data-mgt\\_en.pdf](http://ec.europa.eu/research/participants/data/ref/h2020/grants_manual/hi/oa_pilot/h2020-hi-oa-data-mgt_en.pdf).
- Mons B, Neylon C, Velterop J, Dumontier M, da Silva Santos LOB, Wilkinson MD. Cloudy, increasingly FAIR; revisiting the FAIR Data guiding principles for the European Open Science Cloud. *Information Services & Use* 2017;37(1):49–56. <http://dx.doi.org/10.3233/ISU-170824>.
- Burgelman JC, Pascu C, Szkuta K, Schomberg RV, Karalopoulos A, Repanas K, et al. Open Science, Open Data, and Open Scholarship: European Policies to Make Science Fit for the Twenty-First Century. *Frontiers in Big Data* 2019 dec;2.
- figshare; <https://figshare.com/>.
- Zenodo; <https://zenodo.org/>.
- Dryad; <https://datadryad.org/>.
- Durinx C, McEntyre J, Appel R, Apweiler R, Barlow M, Blomberg N, et al. Identifying ELIXIR Core Data Resources [version 2; referees: 2 approved]. *F1000Research* 2017;5(2422). <http://dx.doi.org/10.12688/f1000research.9656.2>.
- Coordinators NR. Database resources of the national center for biotechnology information. *Nucleic acids research* 2018;46(Database issue):D8.
- Leinonen R, Sugawara H, Shumway M, Collaboration INSD. The sequence read archive. *Nucleic acids research* 2010;39(suppl\_1):D19–D21.
- Submission System of the European Nucleotide Archive; <https://www.ebi.ac.uk/ena/submit>.
- Submission System of the European Variation Archive; <https://www.ebi.ac.uk/eva/?Submit-Data>.
- Krajewski P, Chen D, Ćwiek H, van Dijk ADJ, Fiorani F, Kersey P, et al. Towards recommendations for metadata and data handling in plant phenotyping. *Journal of Experimental Botany* 2015;66(18):5417–5427. <http://dx.doi.org/10.1093/jxb/erv271>.
- ISA-Tab-for-plant-phenotyping; <https://github.com/MIAPE/ISA-Tab-for-plant-phenotyping>.
- Arend D, Lange M, Chen J, Colmsee C, Flemming S, Hecht D, et al. e!DAL – a framework to store, share and publish research data. *BMC Bioinformatics* 2014 Jun;15(1):214. <http://dx.doi.org/10.1186/1471-2105-15-214>.
- Brase J. DataCite—A global registration agency for research data. In: 2009 Fourth International Conference on Cooperation and Promotion of Information Resources in Science and Technology IEEE; 2009. p. 257–261.
- Arend D, Junker A, Scholz U, Schüler D, Wylie J, Lange M. PGP repository: a plant phenomics and genomics data publication infrastructure. *Database* 2016;2016:baw033. <http://dx.doi.org/10.1093/database/baw033>.
- Halewood M, Chiurugwi T, Sackville Hamilton R, Kurtz B, Marden E, Welch E, et al. Plant genetic resources for food and agriculture: opportunities and challenges emerging from the science and information technology revolution. *New Phytologist* 2018;217(4):1407–1419.
- Milner SG, Jost M, Taketa S, Mazón ER, Himmelbach A, Oppermann M, et al. Genebank genomics highlights the diversity of a global barley collection. *Nature genetics* 2019;51(2):319.
- Rosenqvist E, Großkinsky DK, Ottosen CO, van de Zedde R. The Phenotyping Dilemma—The Challenges of a Diversified Phenotyping Community. *Frontiers in Plant Science* 2019;10:163.
- Oppermann M, Weise S, Dittmann C, Knüpffer H. GBIS: the information system of the German Genebank. *Database* 2015;2015.
- Schmutzer T, Bolger ME, Rudd S, Chen J, Gundlach H, Arend D, et al. Bioinformatics in the plant genomic and phenomic domain: The German contribution to resources, services and perspectives. *Journal of Biotechnology* 2017;261:37 – 45.
- Bastow R, Leonelli S. Sustainable digital infrastructure. *EMBO reports* 2010 sep;11(10):730–734.
- Leinonen R, Akhtar R, Birney E, Bower L, Cerdeno-Tárraga A, Cheng Y, et al. The European Nucleotide Archive. *Nucleic Acids Research* 2011;39:D28–D31. <http://dx.doi.org/10.1093/nar/gkq967>.
- Apweiler R, Bairoch A, Wu CH, Barker WC, Boeckmann B, Ferro S, et al. UniProt: the Universal Protein knowledge-base. *Nucleic Acids Research* 2004;32:D115–D119. <http://dx.doi.org/10.1093/nar/gkh131>.
- Perez-Riverol Y, Csordas A, Bai J, Bernal-Llinares M, Hewapathirana S, Kundu DJ, et al. The PRIDE database and related tools and resources in 2019: improving support for quantification data. *Nucleic Acids Research* 2018 nov;47(D1):D442–D450.
- Le Novère N, Bornstein B, Broicher A, Courtot M, Donizelli M, Dharuri H, et al. BioModels Database: a free, centralized database of curated, published, quantitative kinetic models of biochemical and cellular systems. *Nucleic Acids Research* 2006;34(suppl 1):D689–D691. <http://dx.doi.org/10.1093/nar/gkj092>.
- Crosswell LC, Thornton JM. ELIXIR: a distributed infrastructure for European biological data. *Trends in Biotechnology* 2012 may;30(5):241–242.
- Schomburg I, Jeske L, Ulbrich M, Placzek S, Chang A, Schomburg D. The BRENDA enzyme information system – From a database to an expert system. *Journal of Biotechnology* 2017;261:194 – 206. <http://dx.doi.org/10.1016/j.jbiotec.2017.04.020>.
- Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Research* 2013;41(D1):D590–D596. <http://dx.doi.org/10.1093/nar/gks1219>.
- Chavan V, Penev L. The data paper: a mechanism to incentivize data publishing in biodiversity science. *BMC Bioinformatics* 2011 Dec;12(15):S2.
- Fenner M, Crosas M, Grethe JS, Kennedy D, Hermjakob H, Rocca-Serra P, et al. A data citation roadmap for scholarly data repositories. *Scientific Data* 2019;6(1):28. <https://doi.org/10.1038/s41597-019-0031-8>.
- Haak LL, Fenner M, Paglione L, Pentz E, Ratner H. ORCID: a system to uniquely identify researchers. *Learned Publishing* 2012;25(4):259–264. <http://dx.doi.org/10.1087/>

- 20120404.
35. Dreyer B, Hagemann-Wilholt S, Vierkant P, Strecker D, Glagla-Dietz S, Summann F, et al. Die Rolle der ORCID iD in der Wissenschaftskommunikation: Der Beitrag des ORCID-Deutschland-Konsortiums und das ORCID-DE-Projekt. *ABI Technik* 2019 jul;39(2):112–121.
  36. Consortium WWW, et al. JSON-LD 1.0: a JSON-based serialization for linked data 2014;<http://hdl.handle.net/10421/7478>.
  37. Guha RV, Brickley D, Macbeth S. Schema.org: evolution of structured data on the web. *Communications of the ACM* 2016 jan;59(2):44–51.
  38. Weibel S. The Dublin Core: A Simple Content Description Model for Electronic Resources. *Bulletin of the American Society for Information Science and Technology* 1997;24(1):9–11. <http://dx.doi.org/10.1002/bult.70>.
  39. Linden M, Prochazka M, Lappalainen I, Bucik D, Vyskocil P, Kuba M, et al. Common ELIXIR Service for Researcher Authentication and Authorisation [version 1; peer review: 3 approved, 1 approved with reservations]. *F1000Research* 2018;7(1199).
  40. Hardt D; Internet Engineering Task Force. The OAuth 2.0 Authorization Framework (RFC 6749) 2012;<http://dx.doi.org/10.17487/RFC6749>.
  41. Tauch A, Al-Dilaimi A. Bioinformatics in Germany: toward a national-level infrastructure. *Briefings in Bioinformatics* 2017 04;20(2):370–374.
  42. Cousijn H, Feeney P, Lowenberg D, Presani E, Simons N. Bringing Citations and Usage Metrics Together to Make Data Count. *Data Science Journal* 2019;18(1).
  43. Konkiel S. Tracking citations and altmetrics for research data: Challenges and opportunities. *Bulletin of the American Society for Information Science and Technology* 2013;39(6):27–32.
  44. Tenopir C, Dalton ED, Allard S, Frame M, Pjesivac I, Birch B, et al. Changes in Data Sharing and Data Reuse Practices and Perceptions among Scientists Worldwide. *PLOS ONE* 2015 aug;10(8):e0134826.
  45. Parsons MA, Duerr RE, Jones MB. The History and Future of Data Citation in Practice. *Data Science Journal* 2019;18.
  46. Data citation needed. *Scientific Data* 2019;6(1):27.
  47. Hahnel M, Valen D. How to (Easily) Extend the FAIRness of Existing Repositories. *Data Intelligence* 2019 nov;p. 192–198.



GigaScience, 2018, 1–??

doi: [xx.xxxx/xxxx](#)

Manuscript in Preparation  
Technical Note

## TECHNICAL NOTE

# The on-premise data sharing infrastructure e!DAL: Foster FAIR data for faster data acquisition

Daniel Arend<sup>1\*</sup>, Patrick König<sup>1</sup>, Astrid Junker<sup>1</sup>, Uwe Scholz<sup>1</sup> and Matthias Lange<sup>1</sup>

<sup>1</sup>Leibniz Institute of Plant Genetics and Crop Plant Research (IPK) Gatersleben

\*arendd@ipk-gatersleben.de

## Abstract

The FAIR data principle as a commitment to support long-term research data management is widely accepted in the scientific community. Although the ELIXIR Core Data Resources and other established infrastructures provide comprehensive and long-term stable services and platforms for FAIR data management, a notifiable number of research data is still hidden or under risk of getting lost. Currently, high-throughput plant genomics and phenomics technologies are producing research data in abundance, the storage of which is not covered by established core databases. Due to its high volume and strong heterogeneity, resulting in missing infrastructures these data sets often remain unpublished. To master these challenges the ELIXIR Germany/de.NBI service *Plant Genomic and Phenomics Research Data Repository (PGP)* implements a “bring the infrastructure to the data” approach which allows to keep the research data in place and wrap them by a FAIR-aware software infrastructure. This article presents new features of the PGP infrastructure as a best practice on how to easily set up FAIR-compliant and intuitive research data services. Furthermore, the integration of the ELIXIR AAI and data discovery services are introduced as means to lower technical barriers and to increase the visibility of research data.

**Key words:** research data management; FAIR principles; digital object identifier; plant genomics and phenomics

## Introduction

The FAIR (Findability, Accessibility, Interoperability, and Reusability) principles, drafted by the FORCE11 workgroup in 2015 [?] and published in 2016 by Wilkinson et al. [?] are widely accepted and are increasingly adopted in the management policies of research data. The scientific community is showing a rising awareness of the scientific value of reusable research data. This has already resulted in FAIR principles being formally accepted in several data management guidelines, e.g. in the Horizon2020 program [?] of the European Commission, and integrated into the research funding policy [?]. Its technical implementation is supported by data repositories, which store and share research data in a FAIR manner. Those can be classified into (i) general purpose data repositories, e.g. figshare [?], Zenodo [?] and Dryad [?], (ii) core data deposition databases, i.e. ELIXIR deposition databases for life science data [?], NCBI database resources [?] and (iii) institutional

databases and repositories. All have in common that the research data has to be transferred by its owner from the place of data generation to these repositories. This involves considerable effort for data compilation, cleansing, homogenisation, metadata enrichment, formatting and upload. As a result, the published datasets are condensed and generally limited to insufficiently documented supplement material for publications in scientific journals. In the case that data is prepared for a submission to database systems, e.g. the EBI and NCBI core data resources, Bioinformaticians are charged and trained. Examples are the preparation of data for submission to the EBI ENA archive [?], the European Variation Archive (EVA) [?] or the preparation of ISA-TAB compatible data submission for plant phenotyping data [?]. Alternatively, project-related or institutional data repositories could be set up. This in turn requires skilled technicians and computer scientists as well as long-term access to appropriate network and storage infrastructure. Such repositories frequently have a short life time, whether

Compiled on: July 6, 2020.

Draft manuscript prepared by the author.

due to staff fluctuation, long-term maintenance costs and resource consumption. Another reason may be that the repository's niche is too specific to attract substantial data volume, which in turn strongly depends on policies and cost-benefit considerations.

Thus, there is a need for an additional class of repositories that support the data sharing for this class of research data by moving the infrastructure to the data. The concept is to apply an on-premise, infrastructure-to-the-data (I2D) principle. The basic idea of the I2D approach is shown in Figure ???. In contrast to the conventional data publication pipelines to journal accepted databases, which usually involves a time-consuming data upload to an external platform and possibly additional costs depending on the required storage space, the underlying e!DAL software [?] encapsulates an existing storage infrastructure by a data publication layer. This layer is a broker to the DataCite [?] data publication service agent and provides an API and a tooling infrastructure for data submission, DOI delivery, reporting and data quality reviewing. This finally enables the assignment of DOIs with a minimal set of metadata to in-house stored data and its approved FAIR referencing by journals or data lookup services.

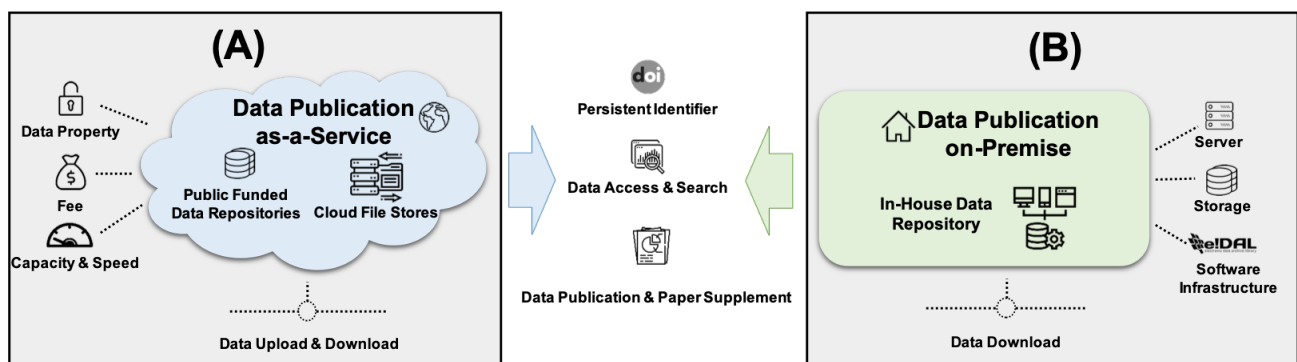
As proof of concept, the Plant Genomics and Phenomics Data Repository (PGP) was implemented [?] to publish digital plant genetic resources (PGR) [?] according to the FAIR principles. Approximately seven million crop accessions are conserved in genbank collections worldwide. The valorisation of PGRs through genotyping and phenotyping is of special focus in the public and private sectors [? ?]. The data management of digital PGRs is identified as one of the most important challenges for a long-term strategy to enhance the productivity, sustainability and resilience of crop varieties and agricultural systems. In contrast to successful studies on genomics-assisted genbank management and the utilization of germplasm collections [?], the special focus of the PGP repository is the publication of buckets of research data that do not fit into classical databases due to their volume, objective, structure or incomplete analysis. Examples are primary data from imaging, field phenotyping, SNP matrices, 3D-plant models, metabolite screenings and environmental sensor data. The experience gained during the four-year operation of the repository has led to a growing acceptance of this approach for the publication of digital PGRs collected in the context of the German Federal ex situ Genebank of Agricultural and Horticultural Crop Species [?]. This experience and the adoption as a service in the European life-sciences Infrastructure for biological Information ELIXIR [?] resulted in novel features, which were implemented with the aim of further improving its acceptance and enabling increased sharing of digital PGRs. After an

update to the state of the art, the new features of the e!DAL data sharing software and its application for the publication of digital PGRs will be explained.

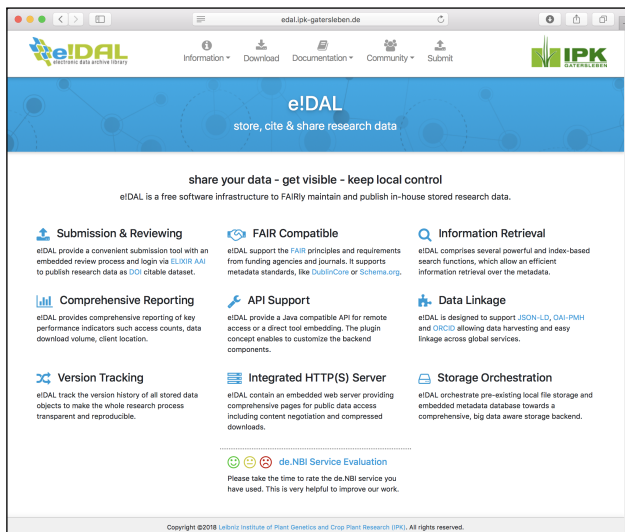
## Related Work

Just like we have many different data types from several domains, there are also a variety of domain-specific archives and information systems. Most of them evolved over many years and they are widely accepted by the research community [?], [ENA](#) for genomic data [?], [UniProt](#) for protein data [?], [PRIDE](#) for proteome data [?], [BioModels](#) for system biology data [?] and many more. Usually also research journals and other publishers request data sharing using these established domain-specific databases. In order to not get lost in the diversity of archives, there are several registries like [re3data.org](#) or [FAIR-sharing.org](#), as well as consortia like [GFBio](#) which collect and categorize repositories to help researchers finding the suitable storage for their data.

Infrastructure programs like the European Open Science Cloud (EOSC) [?], and the European life-sciences Infrastructure for biological Information (ELIXIR) [?] coordinate maintenance and interoperability of research data repositories as federated services by member organisations and hosting institutions. Furthermore, the ELIXIR organisation has the aim to establish a stable and sustainable infrastructure for biological information. In doing so, they define important core resources and tools as a support of the research community [?] like [BRENDA](#) [?] or [SILVA](#) [?]. Most of these systems accept only very specific datasets and require specialised metadata based on schemes that have been improved by the community over years. Unfortunately, there are several, mostly relatively new data types, e.g. plant phenotypic data, which currently do not fit into existing databases, mainly because of their strong heterogeneity and high volume. Public data sharing services like [figshare](#) or [DRYAD](#) provide an alternative solution for publishing these datasets. They are easy to use and have a comprehensive functionality like supporting version controlling and the assigning of persistent identifiers. One important deficiency of such services is the limited free space, which is usually enough for sharing some reduced graphics or aggregated tables, but not for storing large datasets. Furthermore, the establishment and configuration of an own in-house infrastructure based on existing software packages like [CKAN](#) or [Dataverse](#) could overcome this shortcoming, but they require quite a lot of technical prerequisites and know-how.



**Figure 1.** Data-Publication-as-a-Service vs. Data-Publication-on-Premise: Both services feature FAIR data publication, but differ in costs. The as-a-service model (A) costs fee, delegation of data property control and face capacity limits in storage and data upload. The on-premise model (B) keep data in-house, but requires the availability of server and storage hardware as well as the installation of the e!DAL software.



**Figure 2.** Screenshot of the e!DAL project website. On the project website a very detailed description on how to use the e!DAL as well as comprehensive examples for users as well as for developers. Furthermore some short video tutorials and webinar recordings are available to lower the initial barriers to dig into the infrastructure

## Infrastructure

To lower the technical barriers and minimize the effort for scientists to archive and share their research data, we developed the generic **e!DAL software infrastructure** [?]. The usual “Data Publication as-a-Service” procedure includes the transfer of selected datasets to external databases and storage infrastructures after data generation and analysis. In this way research data can be referenced in a future research publication, as show on the left side of Figure ?? (A). In contrast the e!DAL infrastructure provide a “Data Publication on-Premise” approach which enables the publication of locally stored high voluminous research data through the assignment of widely accepted and long-term stable Digital Object Identifiers (DOIs). This is illustrated on the right side of Figure ?? (B). Using DOIs for referencing provides multiple advantages for sharing and accessing research data. Beside adding them as supplements to a research article they can also be the basis for a comprehensive data paper [?]. Furthermore, the well connected infrastructure of the DataCite consortia strongly increase the visibility of the research data assigned with a DOI. It is automatically linked with the ORCID account of the authors, can be found via the DataCite Search and other common search engines or can be harvested via OAI-PMH interface.

Based on the e!DAL infrastructure, the **Leibniz Institute of Plant Genetics and Crop Plant Research (IPK) Gatersleben** and the **German Plant Phenotyping Network (DPPN)** jointly have initiated the **Plant Genomics and Phenomics Research Data Repository (PGP)** [?] as a powerful infrastructure for the publication of comprehensive plant genomics and phenomics research data. The repository covers in particular cross-domain datasets, which are not being published in public repositories for reasons of data volume or data domain, such as phenotyping images, genotyping data, visualizations of morphological models, data from mass spectrometry as well as software and related documents. Doing so, PGP currently provides 200 data records, which can be referenced via DOIs and are annotated with technical metadata. These records comprise more than 1.4 million files with an overall volume of over 2.6 terabytes (see Figure ??). To ensure data discoverability, PGP provides landing pages with JSON-LD formatted metadata and is therefore discoverable through data web crawler services, which follow

the schema.org recommendations, such as Google, Microsoft, Yandex etc. Furthermore, e!DAL implements the OAI-PMH Protocol for Metadata Harvesting from Open Archives Initiative. To support scientists to disseminate their research data the PGP infrastructure is accepted as institutional repository for the **Scientific Data** (Nature Publishing Group) and the **GigaScience** (Oxford Academic) Journals and is registered in [re3data.org](https://re3data.org), [FAIRsharing.org](https://fairsharing.org), [OpenAIRE](https://openaire.org) and [DataCite](https://datacite.org).

The benefits of this wide support of data discovery enabling technologies and data publication in general is proven by the steadily increasing number of dataset accesses. By June 2020, PGP delivered 300 terabytes of data and the provided datasets have been accessed by 100,000 unique clients.

## Improvements

The following section will sum up the main enhancements and updates of the e!DAL infrastructure which comprise new general features, comprehensive changes of several frontend components and important performance improvements. Furthermore, an extensive update due to the latest changes in the Java programming languages and an improved build and deployment process is described.

## Performance

After releasing the first productive version of the PGP repository in 2015, we received many diverse data submissions from several research domains and with very heterogeneous data files. Since then we recognized that the e!DAL infrastructure software scales very well and is able to handle millions of data files, which confirms previous calculations and performance tests [?]. But furthermore it became apparent that sometimes the performance decreases, e.g. for uploading comprehensive datasets with several hundred thousands of small files. Since this is a very common case, e.g. for plant phenotyping datasets, an improvement of the implementation of the e!DAL infrastructure was necessary.

One important feature of e!DAL is the automatic calculation of several essential technical metadata, like the MIME-Type, the data volume or the checksum of every file when storing new datasets. This is convenient, because the user does not need to provide these information on his/her own, but of course these computations are a resource- and runtime-intensive process. Therefore, the functionality to determine the previously mentioned technical metadata and the procedure to transfer the actual binary data have been improved towards a parallel processing of multiple files. This results in a better performance especially on today’s multi-core systems. Furthermore, we optimized several settings for the streaming buffer size and the remote transfer to improve the memory usage and the upload performance for the case of numerous small files. Additionally the checksum calculation was updated to use the more collision resistant SHA-256 algorithm, instead of the older and insecure MD5 function.

## New Features

The previous version of the e!DAL infrastructure already fulfilled several recommendations of the FAIR data principles, such as the support of standardized metadata based on the DublinCore schema or the provision of persistent DOIs for accessing and referencing of research datasets. The e!DAL infrastructure has been further updated to optimize the usability and the general user experience. Additional features were implemented to increase the visibility of published data and the



acceptance of the infrastructure, which in the end also led to be even more FAIR compliant. Thereby the roadmap for scholarly data repositories [?] was taken into account. Subsequently, the most important extensions are described.

### ORCID

To efficiently find and access specific research data files across millions of datasets, persistent identifiers like DOIs or URNs are very helpful and well established. But nevertheless, the research community is also quite large, and sometimes it is very difficult to distinguish data authors because of similar names or to identify the same researcher after he changes his affiliation. With the *Open Researcher and Contributor ID (ORCID)* there is an easy and persistent solution to uniquely identify authors and to solve issues with name ambiguity [?]. An important advantage is the interdisciplinarity, because ORCID is used across nearly all research domains and organizations, e.g. at the mid of 2019 in Germany there were already 150.000 ORCID's registered [?]. By linking authors with publications, affiliations or funding agencies, it helps to find relationships between researchers and their work and the corresponding research data. Since the e!DAL infrastructure is generic and suitable for different kinds of research data, the ORCID system gives us an ideal solution to identify authors and improve the collected metadata for published datasets. Furthermore, the authors and their research data will get a better visibility, due to the connection between the ORCID infrastructure and infrastructure of the DataCite consortium, which is handling the DOIs.

To add the possibility for assigning an ORCID to every data creator or contributor in the e!DAL infrastructure, the PERSON data type in the metadata schema was extended. e!DAL uses the REST API of the ORCID registry to provide the possibility to search for the ORCID of a given name. In addition, it can be validated if an entered ORCID belongs to the corresponding name to prevent an accidentally linking with a wrong ORCID. All these API functions were integrated into the graphical user interface of the data submission tool for the PGP repository. Furthermore, the content pages of published and DOI linked datasets were improved to provide direct links to the ORCID profiles of the associated authors and contributors of the data.

### JSON-LD & DC meta tags

Another method of making research data interoperable as well as machine-readable is to embed the describing metadata using *JavaScript Object Notation for Linked Data (JSON-LD)* format. This approach provides comprehensive possibilities to harvest

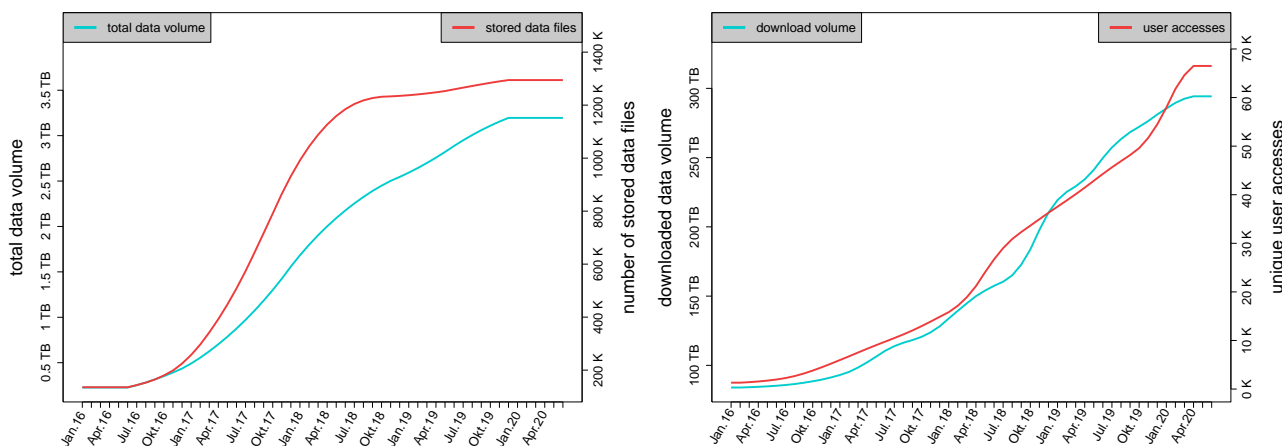
and reuse research data. JSON-LD is a data serialization and exchange method and was developed to be easily embeddable into various systems for providing interoperable web services [?]. The dynamic HTML templates for the content pages of the embedded webserver of e!DAL, which provides the URLs for resolving the assigned DOIs, have been extended accordingly.

```
<script type="application/ld+json">
{
  "@context": "http://schema.org",
  "@type": "Dataset",
  "@id": "https://doi.org/10.5447/IPK/2016/7",
  "name": "Raw images files from quantitative monitoring of...",
  "publisher": {
    "@type": "Organization",
    "name": "IPK Gatersleben"
  },
  "description": "This dataset contains 30426 raw image files...",
  "keywords": "high throughput plant phenotyping, growth protocol...",
  "inLanguage": "en",
  "author": [
    {
      "@type": "Person",
      "givenName": "Astrid",
      "familyName": "Junker",
      "address": "IPK Gatersleben"
    },
    {
      "@type": "Person",
      "givenName": "Thomas",
      "familyName": "Altmann",
      "address": "IPK Gatersleben"
    }
  ],
  "contributor": [
    {
      "@type": "Person",
      "givenName": "Thomas",
      "familyName": "Altmann",
      "address": "IPK Gatersleben"
    }
  ]
}
```

**Listing 1.** Reduced example of the JSON-LD data from the content page of a DOI assigned with e!DAL, which is stored in the PGP repository

Listing ?? show an example for the JSON-LD description of a dataset in the PGP repository. The attributes are based on the [schema.org](http://schema.org) ontology, which is a well-established and community-driven vocabulary used to structure digital data on websites. It is used and harvested by several common search engines [?] and provides an interoperability between dataset from separated resources and platforms.

Another alternative to JSON-LD are so-called HTML meta tags. They are embedded in the <head> section of an HTML-document and also allow to harvest the metadata and describe connections between datasets from different infrastructures. As the metadata schema of the e!DAL infrastructure is already inspired by the DublinCore metadata schema [?] the embedded HTML templates for the content pages of published datasets were extended to provide the technical metadata of every object also as HTML meta tags (see Listing ??).



**Figure 3.** Access numbers and stock overview of the eDAL-PGP repository since 2016. The left diagram shows the development of the data volume and number of stored files, which were provided by the repository. The diagram on the right side the constantly increasing number of accesses and downloads of published datasets is shown.

```

<meta name="DC.Title" content="Screening_of_wild_potato_genetic...">
<meta name="DC.Identifier" content="https://doi.org/10.5447/IPK/2019/1">
<meta name="DC.Publisher" content="e!DAL_-_Plant_Genomics_and_Phenomics...">
<meta name="DC.Language" content="en">
<meta name="DC.Description" content="This_data_set_contains_results_of...">
<meta name="DC.Rights" content="CC_BY-NC-SA_4.0">
<meta name="DC.Creator" content="Bachmann-Pfabe,_Silvia...">
<meta name="DC.Contributor" content="Dehmer,_Klaus...">
<meta name="DC.Subject" content="Phytophthora_infestans">
<meta name="DC.Subject" content="germplasm_collection">

```

**Listing 2.** Reduced example of the DublinCore Meta-Tags from the content page of a DOI assigned with e!DAL

### Content Negotiation

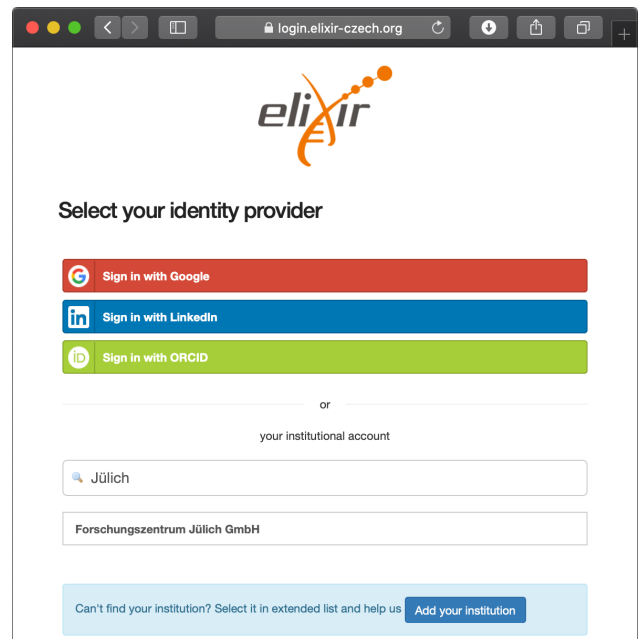
Persistent DOIs provide a solution for long-term stable resolvability and referencing of all published datasets. In addition, for several reasons such as citing the datasets or harvesting the metadata, it is necessary to provide content negotiation to serve resources in different formats. Therefore the possibility to get different representations of the public datasets stored in an e!DAL infrastructure was implemented and can be used by several export functions, which were added on the corresponding content pages as shown in Figure ???. They provide the option to get textual representations, citation formats like BibTeX or RIS and linked data formats like schema.org/JSON-LD and RDF for every dataset. Due to the fact that the [DataCite service](#) already provides a content negotiation feature, it was not necessary to implement a separate function for the embedded webserver of e!DAL. Instead, the HTTP handler uses the provided function for the different formats via a REST call and redirects the responses to the e!DAL infrastructure.



**Figure 4.** Screenshot of an example datasets in the PGP repository. The red colored rectangle mark the e!DAL embedded functions for the content negotiation.

### Elixir AAI

The e!DAL infrastructure provides a flexible and embedded security concept based on the Java Authentication and Authorization Service (JAAS). To provide the research data management and publication capabilities to a wide range of users from universities, research institutes or further organisations, a new login module using the ELIXIR Authentication and Authorization Infrastructure (AAI) [?] was implemented. It was designed to provide a single sign-on service for authenticating researchers to services, which are a part of the ELIXIR portfolio. Doing so, it combines the huge amount of existing organisational identity providers from institutes that are associated with ELIXIR under one roof.



**Figure 5.** Integration of the ELIXIR AAI login dialog into the e!DAL infrastructure. The ELIXIR AAI provides beside a large collection of institutional identity providers, which can be easily found via the small search box, also several third party login opportunities based on Google, ORCID or LinkedIn. This provides the users a very flexible authentication infrastructure.

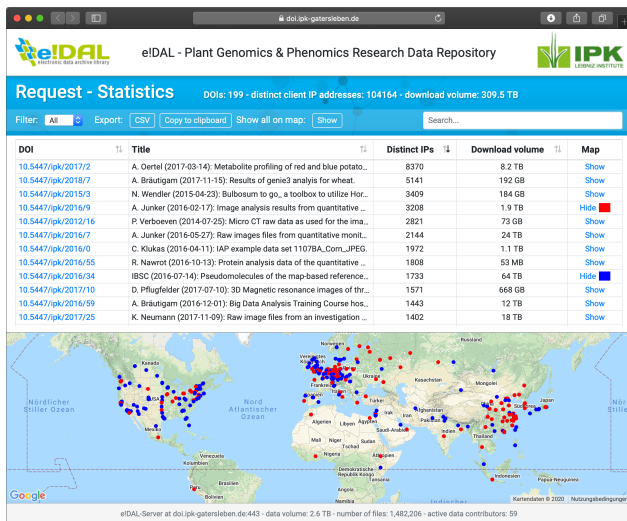
The new e!DAL login module follows the OAuth protocol [?] to authenticate users over the ELIXIR AAI and automatically receive their email address, which is necessary for the communication between the data submitting researcher and the reviewers of the embedded review process. As the first use case, the new ELIXIR AAI based login was integrated into the PGP repository to open the infrastructure and the data submission process for offering the service to a wide range of researchers without the need of creating a separate account. The ELIXIR AAI allows researchers to use their existing organisational accounts (see Figure ??), which lowers the barrier to use the infrastructure and to reach a larger group of data providers.

Furthermore, with the opportunity to use the ELIXIR AAI, the already low effort, which is necessary to establish further e!DAL installations, was reduced. Therefore at the end of 2018 a further e!DAL based repository at the [Jülich Plant Phenotyping Center \(JPPC\)](#) was established using the ELIXIR AAI login provider.

### Amended frontend

The [Apache Velocity](#) template engine is used to render all HTML-based content of the e!DAL embedded webserver like the landing pages of published datasets and e-mail messages. This prevents the infrastructure from storing a massive amount of very similar websites and text drafts, which saves storage and provides a high performance for delivering content via the HTTP handler. All websites are provided dynamically on demand and created from only a few reusable templates.

For the latest e!DAL version all content pages and the underlying templates were fully redesigned to provide a pleasing visual look and functional user experience. By using frontend frameworks and libraries like [Bootstrap](#) and [jQuery](#) it is ensured that the user interface is responsive and working on both modern desktop browsers as well as on mobile devices. Figure ?? shows the new layout as an example screenshot of the embedded access statistic page of the PGP repository. Together with the new design for the frontend components of e!DAL, also the project website was renewed to provide comprehensive



**Figure 6.** Example screenshot of the e!DAL embedded report page for the e!DAL-PGP repository showing the new layout and user interface. Several simple filter and search possibilities help to look for specific datasets. An integrated world map gives the user an idea on the dissemination of the data.

information for the user and for developers in a more concise manner.

## Deployment and Usability

Since the last major release of the e!DAL infrastructure software a lot of optimizations and several new functionalities, which were described in the previous sections, have been implemented. Together with these improvements, changes in the general build and release process and in the usability have also been integrated. The most relevant of them are explained subsequently.

### Gradle Multi-Build Project

After using the [Maven](#) build system Build System for several years for developing and releasing the e!DAL software components a change to the [Gradle](#) build tool was performed. Due to the constant increasing size of the project and the source code, because of new functionalities, several extensions and additional unit tests to guarantee a high software quality the build process using Maven takes quite a long time. This makes the regular release of stable versions very time-intensive. Furthermore, the build configuration became more complex and difficult to maintain. Gradle is strongly focused on a fast and specific build cycle. It supports multi-core systems to a high degree and allows e.g. the execution of several test suites in parallel. With the change to the build infrastructure, we also decided to redesign the entire project build hierarchy and created a multi-build project for the e!DAL infrastructure. It contains the main API components including the reference implementation as well as the components for the server-client architecture, which is directly based on this core implementation. This approach massively accelerates the build time, simplifies the maintenance and allows a more frequent deployment of new versions. The project is now available in a new [BitBucket repository](#).

Nevertheless, the API is still released as an artifact in the central Maven Repository and can be integrated into other software projects using Maven or Gradle, as shown in Listing ??.

```
----- Maven - 'pom.xml' -----
<project>
<dependencies>
<dependency>
<groupId>de.ipk-gatersleben</groupId>
```

```
<artifactId>eDAL-MetaDataAPI | eDAL-MetaDataAPI-Server | eDAL-
MetaDataAPI-Client</artifactId>
<version>3.0.2</version>
</dependency>
<dependencies>
</project>
```

```
----- Gradle - 'build.gradle' -----
repositories {
mavenCentral()
}
dependencies {
compile 'de.ipk-gatersleben:eDAL-MetaDataAPI:3.0.2'
compile 'de.ipk-gatersleben:eDAL-MetaDataAPI-Server:3.0.2'
compile 'de.ipk-gatersleben:eDAL-MetaDataAPI-Client:3.0.2'
}
```

**Listing 3.** Integration of the e!DAL components into the configuration of an Apache Maven or Gradle based project

### OS specific executables

Due to the complete new development and release cycle by Oracle, the [Java](#) programming environment, which is the basis for the e!DAL infrastructure, changes a lot in recent years. In addition, the comprehensive redesign and reconstruction of the language itself, like the introduction of the new module concept or the removal of popular and formerly native APIs and frameworks like [JavaFX](#) or the [Java Network Launching Protocol \(JNLP\)](#) were some very substantial changes. This strongly influences the e!DAL implementation, because they were also a significant part of the previous version. Unfortunately this impeded at some points the further development of e!DAL infrastructure, because a lot of the used frameworks and libraries needed several months to update their code to be compatible with the latest Java versions. With the new version 3.0.0 the e!DAL infrastructure is fully based on the [Java Runtime Environment \(JRE\) 12](#). Therefore some comprehensive changes were necessary. In order to run e!DAL with the different existing runtimes, e.g. the official runtime from Oracle, but also the alternative and widely used [OpenJDK](#), it was necessary to integrate the [JavaFX](#) library directly into the implementation. This increases the actual size of the API package, but it makes the infrastructure much more compatible and even more independent from the system preconditions than before.

The removal of the support for the popular and well known [JNLP](#) was also a high challenge, because the [Java webstart](#) tool was used to give the user an intuitive and platform-independent way to run the graphical data submission tool. Nevertheless this solution also provides some shortcomings like the need to provide an installed and compatible Java runtime. With the recently developed [jpackage](#) Java provides a powerful tool to pack self-contained applications along with a suitable JRE. We used [jpackage](#) to create a full image of the e!DAL data submission tool together with a reduced JRE, which contains only the necessary java modules and provides separate executables for the most common operation system (Windows, Unix, MacOS). This provides a very convenient usability for data submitter and makes the infrastructure again more compatible and independent from the given system preconditions of the users.

### Web-based submission application

In parallel to the update process due to the previously mentioned changes in the [Java Runtime Environment](#) and the development of the build process to create the self-executable applications for the submission dialog, a new web-based application was implemented to provide an alternative opportunity to upload research data to an e!DAL based infrastructure. The goal was the deployment of a user-friendly web application with the similar functionality of the corresponding desktop tool, but without the need to download the application as an executable or additional plugins. The [Vaadin](#) framework for Rich Internet Applications (RIA) was used for the implementation. Figure ??

**Plant Genomics & Phenomics Research Data Repository**  
Data Submission (BETA)

**Title**  
Image analysis results from quantitative monitoring of 484 *Arabidopsis thaliana* plants

**Description**  
This dataset contains image analysis results based on 30426 raw image files from a high-throughput plant phenotyping experiment. Here, the growth and development of 484 *Arabidopsis thaliana* plants (C24 genotype) in respect to plant movement and soil cover has been monitored over 55 days. This dataset contains 29724 feature-annotated result images. The features are plant height and volume.

First Name	Last Name	ORCID	Legal Name	Address	ZIP	Country	Type	Remove
Daniel	Arend	0000-0002-2455-5938	Leibniz Institute of Plant Genetics and Crop Plant Research		06466	Germany	Creator	
Astrid	Junker	0000-0002-4656-0308	Leibniz Institute of Plant Genetics and Crop Plant Research		06466	Germany	Creator	

**License**  
 CC0 1.0 Universal (Creative Commons Public Domain Dedication)  
 CC BY 4.0 (Creative Commons Attribution)  
 CC BY-SA 4.0 (Creative Commons Attribution-ShareAlike)  
 CC BY-ND 4.0 (Creative Commons Attribution-NoDerivatives)  
 CC BY-NC 4.0 (Creative Commons Attribution-Non-Commercial)  
 CC BY-NC-SA 4.0 (Creative Commons Attribution-Non-Commercial-ShareAlike)  
 CC BY-NC-ND 4.0 (Creative Commons Attribution-Non-Commercial-NoDerivatives)

**Embargo Date**

**Next**

**Figure 7.** Example screenshot of the web-based submission application for the e!DAL-PGP repository. The form like user interface is very intuitive and similar to the desktop application. It helps the user to fill in the needed metadata by executing several checks to guarantee the completeness of the metadata.

shows a screenshot of the web application. By using several REST APIs, e.g. from the ORCID Registry or the ELIXIR AAI, a light-weight application could be created providing the same functionality as the full desktop client. Furthermore, users now have the possibility to submit research data also from mobile devices or other browser compatible devices. The only small shortcoming of the data submission via the web application is currently that not all browsers support the upload of comprehensive file folders. The latter is only possible if a recent version of Google Chrome or Mozilla Firefox is used. Other web browsers only allow the upload of single files.

## Results

In this article the basic overall ‘on-premise’ data management and publication concept of the e!DAL infrastructure as well as several new features and technical developments were presented. As a result, e!DAL matured to a comprehensive and FAIR-compliant infrastructure, while always keeping the focus on simple and flexible setup and integration into existing infrastructures and into the daily research process. With the described ‘bring the infrastructure to the data’ approach, it differs fundamentally from generic publication platforms like figshare or DRYAD, which can produce, depending on the needed storage, considerable financial costs and time costs for transferring the data. e!DAL allows the usage of available in-house storage capacities, without the need of complex requirements and technical infrastructures or comprehensive adaptations. All functionalities are already included and the provided reference implementation contains required components, such as a database or a webserver. This is a crucial advantage in comparison to other similar software infrastructures, like DataVerse or CKAN, and lowers the barrier to establish a publication infrastructure even for small-size research institutions with limited possibilities and know-how.

Thereby the FAIR compliance can be fulfilled by several e!DAL functions and components:

- **Findable:**

By providing embedded and machine-readable metadata based on standardized established formats, the e!DAL published datasets can be easily found using common search engines like Google or the DataCite Metadata Search. Due to

the widely established and used DOIs, the DataCite consortium is also involved in several projects and interacts with different systems like ORCID, CrossRef or Scholix. This further improves the findability of e!DAL datasets.

- **Accessible:**

The usage of DOIs as persistent identifiers guarantees a long-term stable availability of published datasets. The DataCite resolver for the DOIs allows simple access to the data and reference datasets, e.g. in a research article or as part of data publication. If the storage location of the underlying data is changing, the corresponding DOI remains stable and allows the uninterrupted access to the data by updating the resource path. Even if the data is no longer available the metadata is still accessible via the DataCite registry.

- **Interoperable:**

To provide interoperable datasets and to allow the aggregation of information about the relationship of datasets from different sources, the e!DAL infrastructure supplies embedded metadata on the content pages of every data object. They are stored using standardized formats and vocabularies like JSON/LD or rather schema.org.

- **Reusable:**

By collecting a standardized set of mainly technical metadata e!DAL guarantees a long-term readability and usability of all published datasets. The schema is inspired by the DublinCore metadata format and meets community established standards. Furthermore, a clear and easy license handling allows to assign a suitable license, which defines by whom and how the data can be used. They are available both on the content page of every data object as well as embedded in the HTML sources.

### e!DAL Usage

Established in 2016, the PGP repository is the first productive repository based on the e!DAL infrastructure and a part of the service portfolio of the GCBN unit (German Crop BioGreenformatics Network) [?] of de.NBI (German Network for Bioinformatics Infrastructure) [?], which is the head of ELIXIR Germany. After more than three years of productive usage, the PGP repository currently shares 200 comprehensive, plant-related research datasets containing mainly genomic and phenomic information, but also metabolic datasets or software components and pipelines. They have a total data volume of 2.6 terabytes consisting of more than 1,4 million files. Most of the datasets are part of a corresponding research paper and allow authors from IPK, but also from other institutes, to improve their manuscripts by enriching them with the underlying research data in a FAIR compliant way. Until June 2020 an overall download volume of around 300 terabytes and over 100,000 distinct user accesses were counted, which shows the high visibility of the provided data and the interest of the research community for this kind of data.

The integration of the ELIXIR AAI into the login mechanism of the PGP Repository is a prime example that shows how established platforms can benefit from the ELIXIR network. The provided services contribute to the increase of the visibility, to overcome the obstacles for the use of available infrastructures and to support FAIR compliant access to research data. The support of the ELIXIR single sign-on service enables collaborators to use the PGP repository as a service to publish their research data. Furthermore, the ELIXIR AAI login is fully integrated into the e!DAL infrastructure software, which allows to set up further FAIR in-house repository instances, following the presented I2D approach. Doing so in June 2018 a second repository based on the developed e!DAL infrastructure was established at the *Forschungszentrum Jülich*. Due to the auto-configuring installation it was possible to run the system and provide the submission and review workflow with only a

little effort in time. The integrated ELIXIR AAI login allows researchers from Jülich to use their existing institutional accounts. The complete infrastructure is hosted and maintained by the [Jülich Plant Phenotyping Center \(JPPC\)](#). The process of establishing further e!DAL-based repositories at the [Julius-Kühn Institute](#) and the [Helmholtz Centre München](#) are currently underway.

## Outlook

In this work, we showed the newly designed I2D concept for FAIR compliant data publication by using in-house storage infrastructures and new features of the e!DAL platform. After several years of operating a productive instance of this infrastructure as the basis for the PGP repository, we recorded high numbers of accesses and downloads. Although researchers have more and more possibilities to share their research data with the community, the incentive to do so is still not high enough for some researchers [? ]. In contrast to the common peer-reviewed publication in journals, it is not so easy to measure the impact of research data itself, because the concept of data citation is still not a common practice [? ], but it becomes more and more important and accepted [? ]. It's not only a cultural problem, but also a technical challenge and therefore an issue of practicability[? ]. One of the first metrics to count data citations was the commercial [Data Citation Index](#). But in the meanwhile some free and community initiated projects like [Make Data Count](#) have been developed. Furthermore, popular journals are starting to demand that authors put their research data as data citations in their common reference list [? ]. This facilitates to measure its impact through a citation index and improves the visibility to readers which in turn increases the general acceptance of research data as valuable scientific assets. In future we will investigate several approaches for counting data citations and getting more credit for publishing research data. We plan to integrate a generic and open-source solution into the e!DAL infrastructure to show users comprehensive information how their data is reused and referenced.

The ORCID provide a widely accepted and used solution to unambiguously identify researchers. The integration within the e!DAL infrastructure is very intuitive and facilitates handling of multiple ORCIDs for comprehensive lists of authors. Besides the identification of persons, it can be also quite challenging to handle the diverse affiliations of research institutes, universities or companies with a focus on different scientific topics. Some authors have multiple affiliations, from time to time organizations may be renamed, the official addressee may change due to infrastructural developments or it may happen that an institute will be closed. The [Research Organization Registry \(ROR\)](#) provides an open and sustainable approach, which is led by the community and supported by popular organizations like DataCite or Dryad. The concept of the ROR identifiers is very similar to the ORCIDs and allows to uniquely identify all kinds of research organizations. Therefore, one of the next functional improvements for the e!DAL infrastructure will be the integration using the provided ROR API. This will cause some changes in the basic data structure, which however will result in a much easier and FAIRer way to handle author affiliations [? ].

## Availability of source code and requirements

- Project name: e!DAL (electronic Data Archive Library)
- Current version: 3.0.2
- Project homepage: <https://edal.ipk-gatersleben.de>
- Source Code Repository: [https://bitbucket.org/ipk\\_bit\\_](https://bitbucket.org/ipk_bit_)

[team/electronicdataarchivelibrary](https://bitbucket.org/ipk_bit_team/electronicdataarchivelibrary)

- Operating system(s): platform independent
- Programming language: JVM based (Java 12+)
- JavaDoc: <https://edal.ipk-gatersleben.de/javadoc>
- Artifact repository: Maven Central (<https://mvnrepository.com/artifact/de.ipk-gatersleben>)
- License: GNU General Public License (GPL) Version 3 (<https://www.gnu.org/licenses/gpl-3.0.html>)

The e!DAL [project website](#) provide comprehensive information, which are relevant for users as well as developers. Beside a full Java documentation, further presentations, videos and several code and usage examples are presented. We changed the licence model of e!DAL to GNU General Public License (GPL) Version 3. This aim at maximising the spread of the e!DAL infrastructure in scientific community to foster FAIR principles of in-house stored data and to enable the incorporation of e!DAL into 3rd party software as well.

## Declarations

### List of abbreviations

AAI: Authentication and Authorization Infrastructure;  
 API: Application programming interface;  
 de.NBI: German Network for Bioinformatics Infrastructure;  
 DOI: Digital Object Identifier;  
 DPPN: German Plant Phenotyping Network;  
 FAIR: Findable, Accessible, Interoperable, Reusable;  
 GCBN: German Crop BioGreenformatics Network;  
 I2D: Infrastructure-to-the-Data (I2D);  
 JNLP: Java Network Launching Protocol;  
 JPPC: Jülich Plant Phenotyping Center;  
 JRE: Java Runtime Environment;  
 JSON-LD: JavaScript Object Notation for Linked Data;  
 ORCID: Open Researcher and Contributor ID;  
 PGP: Plant Genomics and Phenomics Research Data Repository;  
 PGR: Plant Genetic Resources;  
 RIA: Rich Internet Application;  
 RMI: Remote Method Invocation;

## Competing Interests

The authors declare no conflict of interest.

## Funding

This work was supported by the German Federal Ministry of Education and Research (BMBF) in frame of the the projects German-Plant-Phenotyping Network - DPPN (FKZ 031A053), Modernste Virtualitäts- und erweiterte Realitäts-Verfahren für den Zyklus von Samen zu Samen - AVATARS (FKZ 031B0770A) and German Network for Bioinformatics Infrastructure - de.NBI (FKZ 031A536A).

## Author's Contributions

Conceptualization: DA,ML  
 Software: DA,PK  
 Investigation: DA  
 Supervision: ML  
 Writing original draft: DA,ML  
 Writing review & editing: All  
 Funding acquisition: US,AJ,ML

## Acknowledgements

We want to thank Thomas Münch, Jens Bauernfeind and Heiko Mieke for the technical support. Furthermore we are very thankful to Benjamin Bruns, Robin Kluth and Michael Seidel for fruitful discussions, testing the infrastructure and formulating feature requests. Last but not least we thank all users and data producers for their valuable feedback.