

# The eDAL Suite: Tools and Concepts for Primary Data Citation

M. Lange, C. Colmsee, S. Flemming, J. Chen, M. Klapperstück, U. Scholz

Research Group Bioinformatics and Information Technology,  
Leibniz Institute of Plant Genetics and Crop Plant Research (IPK), Corrensstraße 3, 06466 Gatersleben, Germany  
Contact: <lange>|<colmsee>@ipk-gatersleben.de

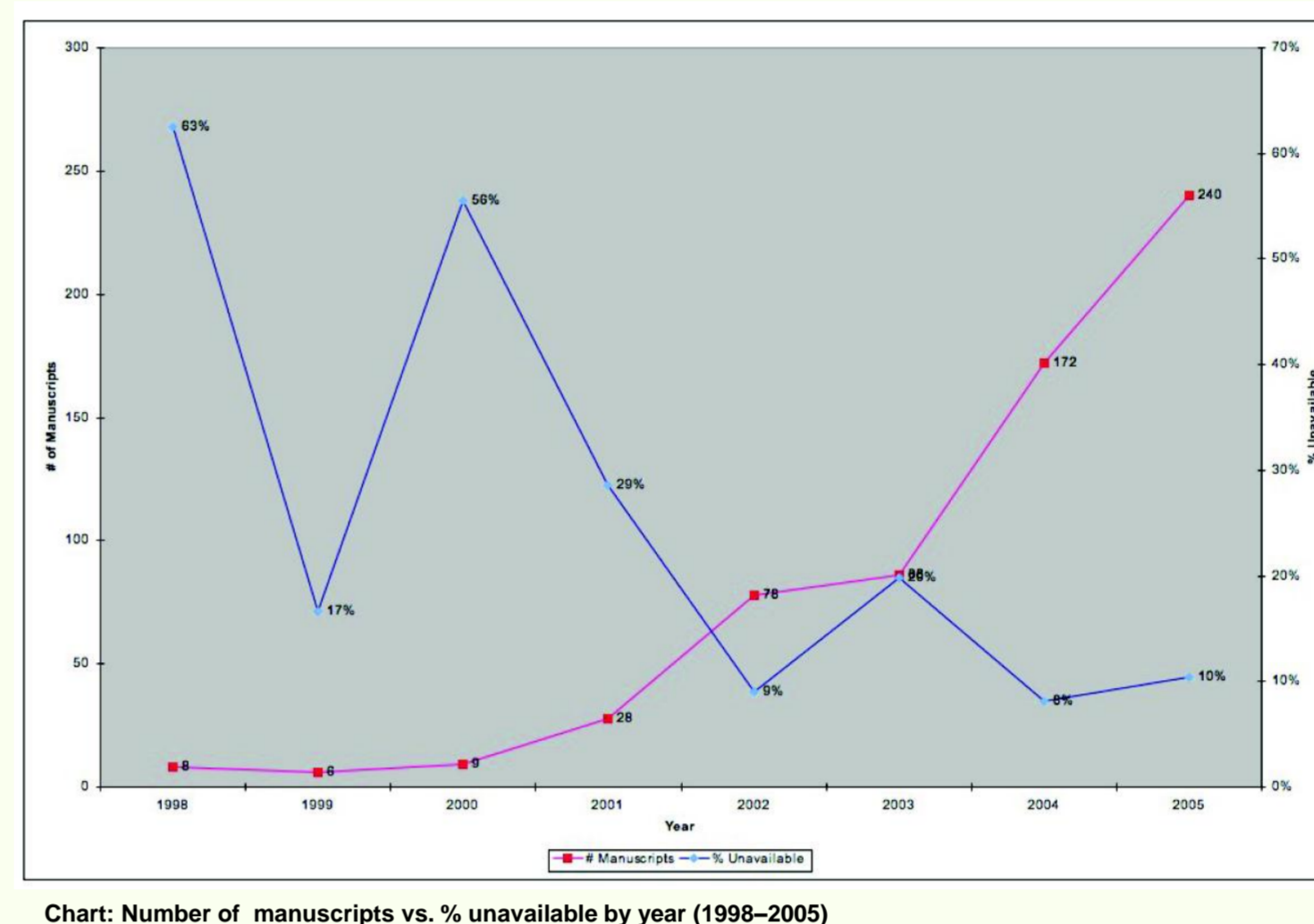


## Motivation

Retrieval and citation of primary data is the important factor in the approaching age of “data science”. Digital data are easily shared, and just as easily wiped or lost. The problem of keeping on-line data accessible and retrievable is especially difficult for SME like plant breeders plant biotech companies as well as research projects in this domain. Intension of eDAL is the provisioning of an information retrieval and data citation infrastructure that meets the requirements of the “data science” age and implements a re-usable platform for data retrieval, data citation, and data publication. Like a shopping cart, the idea is to combine a search engine and a data cart, which retrieves, rank and collect query relevant data from crop plant data centers.

### Distributed, Redundant, Alternative Data @ IPK: 916 data sets for „phosphate ransporter“ in 4 IPK databases

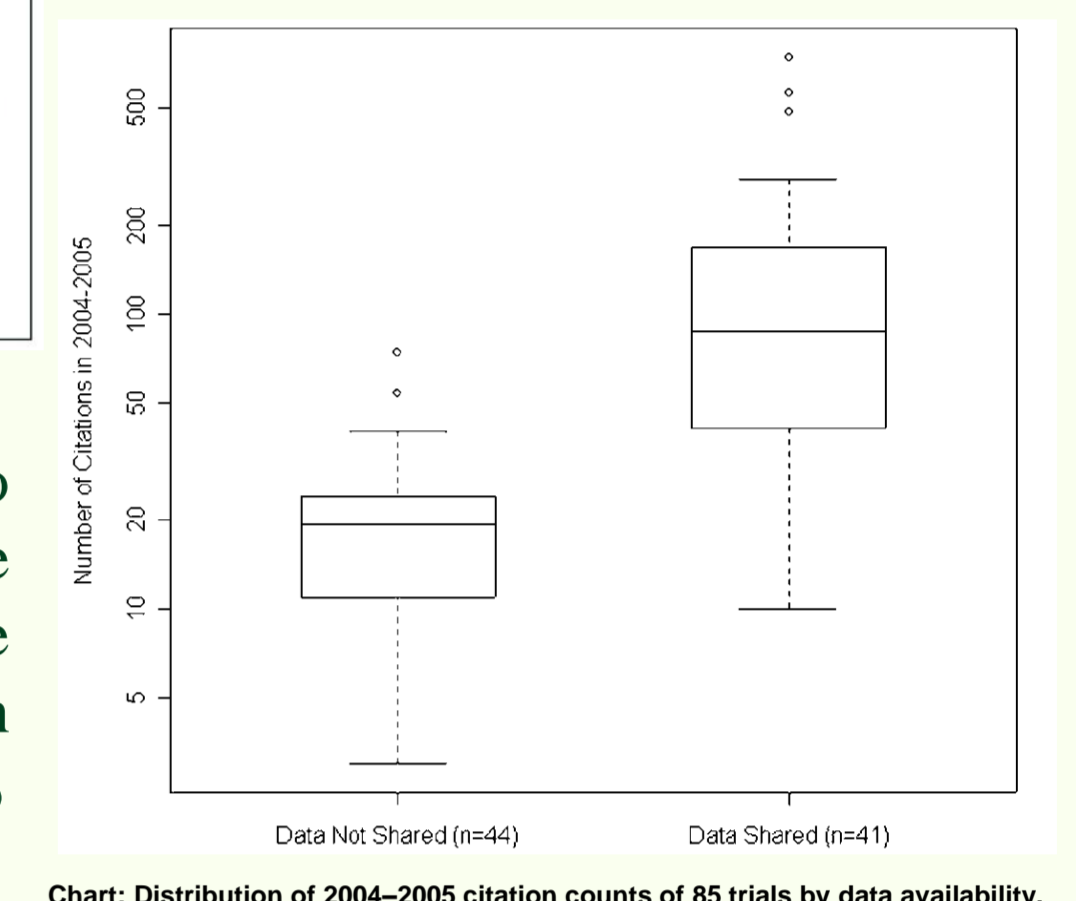
### Availability of Data



The increasing number of digital data should lead into a long-term storage system, within the data is available at all time. But Anderson et al. illustrates the problematic of the loss of digital data within biomedical publications. The reported up to 30% unavailability of supplementary data in 9 years.

### Data Citation Rate

The importance of sharing data within publications is shown in Piwowar et al. for microarray data. Publications comprising shared data have a 69% higher citation rate than publications without shared data.



## State of the Art

Data Citation	Meta Data	Data Access	Search Engine
<p>Proprietary ID's: e.g. Accessions, PubMed-LinkOut</p> <p>Standard ID's: e.g. Life Science Data Identifier (LSID)</p>	<p>Technical meta data: e.g. ISO 15836 (DCMES)</p> <p>Semantic meta data: e.g. EBI Ontology Lookup Service</p>	<p>GBIF/Bio-Moby Data Services</p>	<p>EBI EB-Eye/NCBI Entrez</p>

## Primary Data Management at the IPK

LimsLight: IPK Primary Data Mangement	eDAL-API: Interface to Store, Manage and Annotate Primary Data	Outlook: The CROP-SHOP approach (submitted as BMBF proposal)
<p>The architecture is divided in three parts:</p> <ol style="list-style-type: none"> <li>The client components to upload, search and download the data,</li> <li>The database to maintain meta information (up to a range of Gigabytes),</li> <li>The raw data storage component to archive primary data (volume above Terabyte range)</li> </ol> <p><b>The LIMS Light Web Interface:</b></p> <ul style="list-style-type: none"> <li>Web application based on ORACLE Application Express (APEX)</li> <li>Management of Projects, Experiments, Worksets and files</li> <li>Using controlled vocabulary (Ontology Lookup Service) and free text information to tag the files</li> <li>Simple and advanced search to find stored files</li> </ul> <p><b>Performance Progression</b></p> <ol style="list-style-type: none"> <li>Downloading files with a speed of ~24 MB/s</li> <li>Uploading files with a speed of ~13 MB/s</li> <li>Meta data maintenance for high number of small files causes low upload and download rates</li> </ol>	<p>The eDAL-API takes care of the central aspects of primary data management:</p> <ul style="list-style-type: none"> <li><b>Storage</b> component to archive the primary data</li> <li><b>Persistent identifier</b> to make primary data citable</li> <li><b>Dublin Core</b> for defined meta data information</li> <li><b>JAAS security</b> for data access control</li> <li><b>Versioning</b> of primary data</li> </ul> <p>With the eDAL-API the user can implement his own Primary Data Management system and could realize the information retrieval by using the life science search engine like LAILAPS (Lange et al. 2010).</p>	<p>Data centers and primary data archives (1) provide access to their databases and information systems by a uniform data citation system. The information retrieval component is featured by a search engine (2) and a data cart (3). The search engine retrieves data, which are ranked for their query and user specific relevance profiles. The data cart module collect global object identifiers resulting from a query. The identifiers are linked to administrative meta-data and will be delivered to the end-user for later data processing (4) in adequate data formats.</p>

## References

- M. LANGE et al. (2010) The LAILAPS Search Engine: Relevance Ranking in Life Science Databases. Journal of Integrative Bioinformatics, 7(2):e110, 2010.
- Anderson et al. (2006) On the persistence of supplementary resources in biomedical publications. BMC Bioinformatics 2006, 7:260. doi:10.1186/1471-2105-7-260.
- Piwowar et al. (2007) Sharing Detailed Research Data Is Associated with Increased Citation Rate. PLoS ONE 2(3):e308. doi:10.1371/journal.pone.0000308.

