

**Document Name**    Actual\_state\_of\_research\_data @\_IST\_Austria

**Author**                Jana Porsche  
**Version**                1.0  
**Modified**              12.11.2012  
**Location**              <http://repository.ist.ac.at/id/eprint/103>

## Table of Contents

Purpose.....	2
Study.....	3
Results.....	4
Summary.....	8
Appendix.....	9
Survey text .....	9
Full list of procedures.....	10
Full list of devices/systems.....	11
List of data formats.....	12

## **Purpose**

This document is created as a part of the project “Repository for Research Data on IST Austria”. It summarises the actual state of research data at IST Austria, based on survey results. It supports the choice of appropriate software, which would best fit the requirements of their users, the researchers.

Currently almost 50% of research groups at IST Austria (11 out of 24) are receiving ERC grants and are therefore committed to make their research data publicly available after publication of the related article. Accordingly, an institutional repository will be built to provide the researchers with a platform, which can be used for publication of data based on research done at ST Austria, to fulfil the requests from funding agencies (ERC, FWF, HFSP).

## **Study**

The actual state of the research data at IST Austria was investigated through an on-line survey, distributed via Google Docs. A link to the survey was provided to all professors, who could answer the survey by themselves, or delegate this within their research group.

When the survey was distributed, there were 19 active research groups at IST Austria and all responded. We were looking for one set of answers for every research group. The answers for all questions were defined as free text (no predefined answers), to receive a possibly broad range of answers.

The full text of the survey can be found in the appendix.

## Results

The results are interpreted according to the purpose of the survey - to support the choice of suitable software for publication of research data.

### 1/ How are research data produced in your research process?

Because of the variety of research disciplines, there is a broad range of different procedures: computer simulation, writing code, microscopy, behavioural observation, electrophysiological measurement, DNA sequencing, genomics, in silico simulation, taking pictures, optical density measurement etc.

A full list of the procedures is available in the appendix.

### 2/ Out of which device/system or in which way (for non-digital data) did you get the research data?

Within the 19 research groups, the research data are won through: PC, gas chromatograph, plate reader, patch clamp amplifier, Real Time PCR machine, video camera, photo camera, microscope, GPS, datalogger, CCD camera, fluorescence reader, online genomic database etc.

A full list of the devices and systems is available in the appendix.

### 3/ What kind of research data does your research group produce?

The answers were encoded to 5 main groups of data to simplify the output.

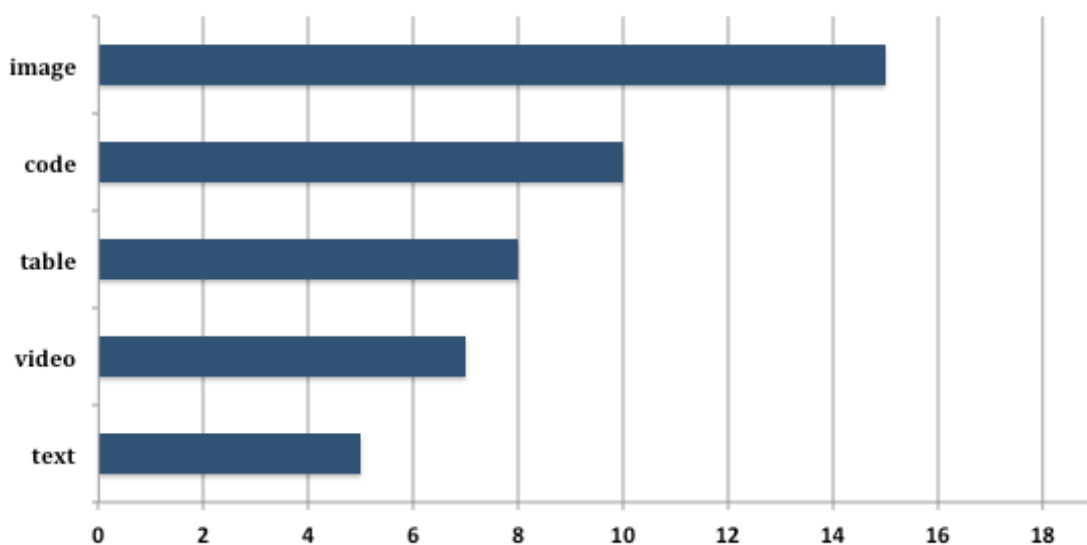
**code** (algorithm, software etc.)

**text** (txt, xml etc.)

**image** (photo, graphic, plot etc.)

**video** (picture stack etc.)

**table** (SPSS, structured textual data, xls)



#### 4/ What formats of research data does your research group produce?

Across the research groups, a very broad range of data formats is produced. To see the frequency of occurrence, the formats are ordered according to usage within the research groups.

##### data format ( number of research groups, using this format)

jpg (9)  
 txt (6)  
 tiff (6)  
 eps (3)  
 pdf (3)  
 tex (3)  
 c (3)  
 cpp (2)  
 doc (2)

ab1 (1), avi (1), bmp (1), cfs (1), clc (1), csv (1), lif (1), ml (1), mov (1), mp4 (1), mpg (1), nb (1), obj (1), ps (1), py (1), r (1), sav (1), spv (1), stk (2), tikz (1), xls (1), xlsx (1), xml (1)

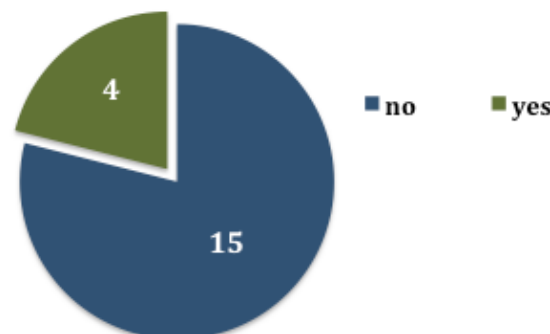
##### Conclusion:

Formats **jpg**, **txt** and **tiff** are the most used and should therefore be tested first, to ensure proper handling within the repository. All mentioned data formats should be generally supported.

Definitions of listed data formats are available in the appendix.

#### 5/ Do you have special requirements for describing the research data?

4 research groups out of 19 would like to add additional description of the data (gene sequences, bacterial strains, brain slices, links to specimens) to the published dataset. 3 of them are from the field of Mathematical and Evolutionary Biology, 1 from Neuroscience.



##### Conclusion:

The repository should enable upload of additional description and materials related to published data.

**6/ What is the approximate amount of research data per month?**

The amount of research data for all 19 research groups together is approximately 1,3 TB in a month.

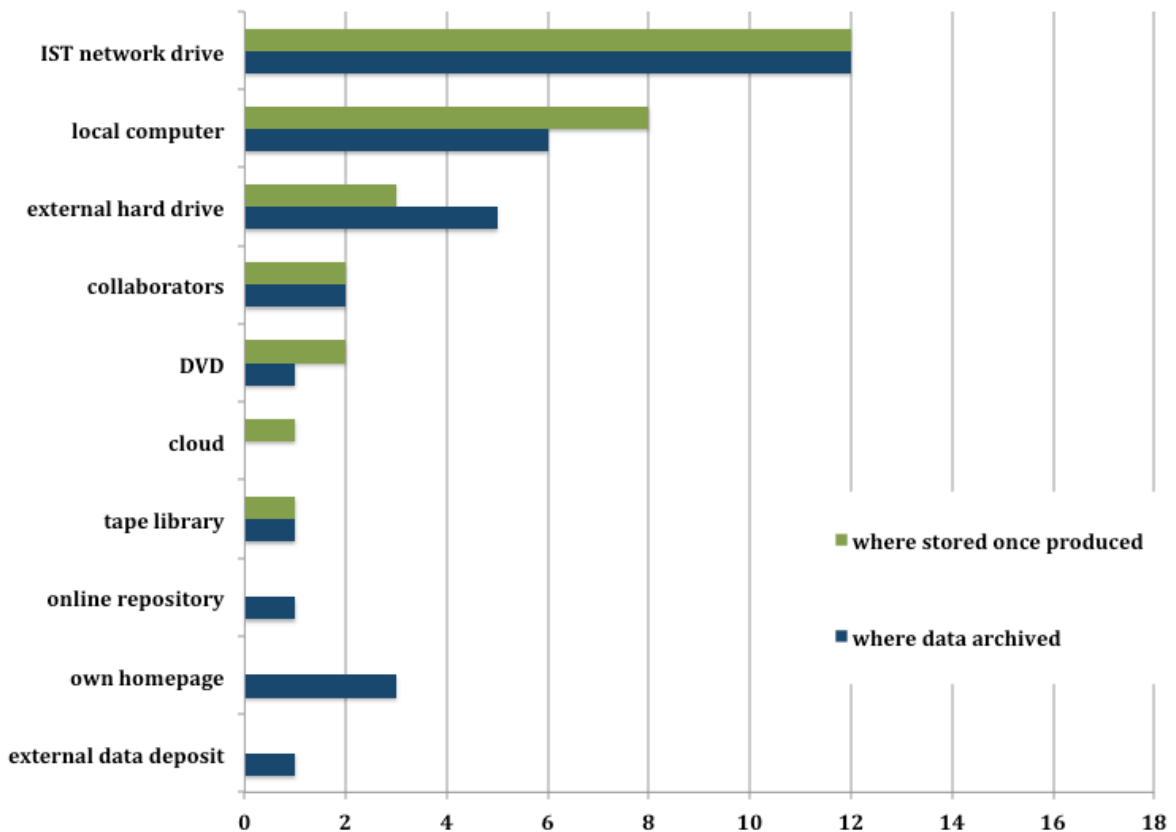
**7/ What is the approximate average data size?**

Approximate average size varies between a few kb to dozens of GB depending on the research field. Neuroscientists are producing hundreds of GB in a month, Computer scientists only a few MB.

**8/ Where do you store (archive) the research data once they were produced?**

**9/ Where do you store (archive) the research data after publication?**

Representation of results for question 8 and 9 is in the joint table below, which shows the differences in strategies.



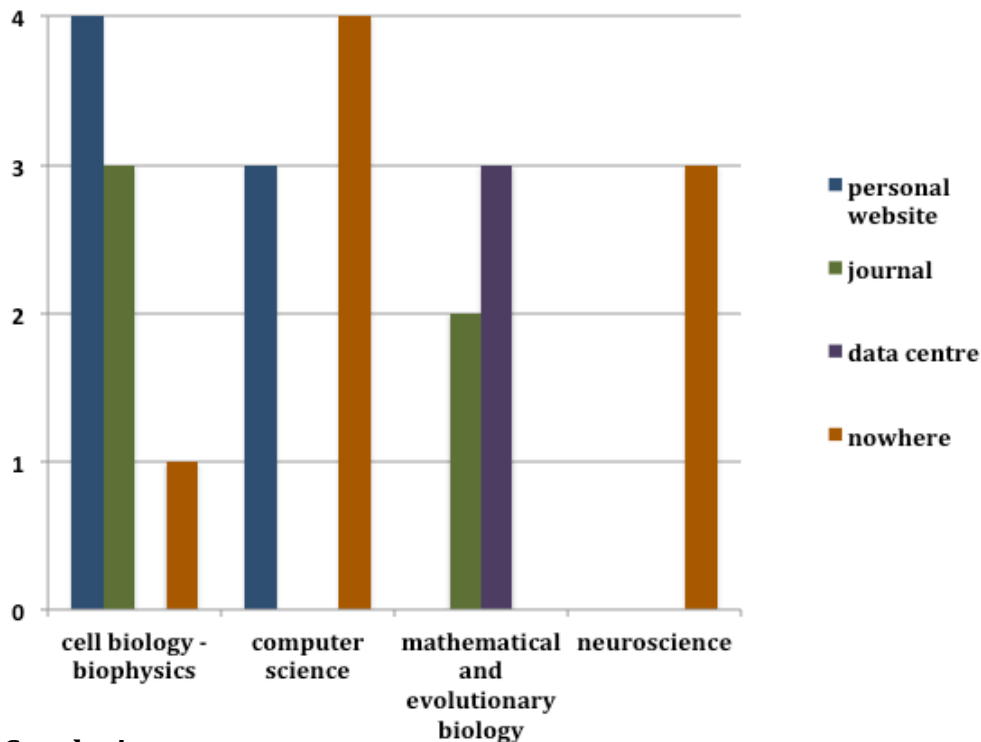
**Conclusion:**

Most of the research groups trust in the IST network drive for temporary data storage and also for data archiving. Still, almost half the research groups use a local computer for data storage, where long-time archiving cannot be ensured. Only a minority of the researches store the data through external methods like data deposits, tape library or cloud.

According to the definition by the DFG<sup>1</sup> in “Safeguarding Good Scientific Practice”<sup>2</sup>: “Primary data as the basis for publications shall be securely stored for ten years in a durable form in the institution of their origin”. This will be better secured by the institutional repository than by individual diverse strategies.

**10/ Where do you upload (publish) your research data after publication?**

The actual experience with publishing data was investigated to see if there is already any experience, or if this is new for the researchers. To simplify the output, the results are grouped by the disciplines to show the difference between the research fields.



**Conclusion:**

Each research field has its own habits for sharing and publishing data, which is visible in the results. In Mathematical and Evolutionary Biology, there is already some experience with publishing research data within journals and also data centres. On the contrary, neuroscientists in our research groups don't publish their data at all.

In Mathematical and Evolutionary Biology, the following data centres are used at present for publishing:

- GenBank**<sup>3</sup> DNA sequence data
- TreeBASE**<sup>4</sup> phylogenetic data
- Dryad**<sup>5</sup> other types of data
- NCEAS Data Repository**<sup>6</sup> other types of data from NCEAS funded activities

<sup>1</sup> Deutsche Forschungsgemeinschaft = German Research Foundation  
<sup>2</sup> [http://www.dfg.de/download/pdf/dfg\\_im\\_profil/reden\\_stellungnahmen/download/empfehlung\\_wiss\\_praxis\\_0198.pdf](http://www.dfg.de/download/pdf/dfg_im_profil/reden_stellungnahmen/download/empfehlung_wiss_praxis_0198.pdf)  
<sup>3</sup> <http://www.ncbi.nlm.nih.gov/genbank/>  
<sup>4</sup> <http://treebase.org>  
<sup>5</sup> <http://datadryad.org/>  
<sup>6</sup> <http://knb.ecoinformatics.org/knb/style/skins/nceas/>

### 11/ Do you upload (publish) all of the research data or only a part of them after publication?

As it is already visible in the answers on the previous question, 8 out of 19 research groups don't publish their research data at the moment. 9 groups publish part of their data (final movie, part of sequence data, supplementary material, rendered video, images) and 2 groups do publish the whole primary data.



### Summary

We started the project „Research Data @ IST Austria“ with an on-line survey to investigate the actual state of research data at our institute. Thanks to 100% response rate, we could work with all individual details and needs through all available research disciplines on IST Austria.

As it is visible from the results, it is important to offer researchers the possibility to publish their research data directly in an institutional repository with the opportunity to save links to published data in well-established subject repositories (data centres). At the moment, the individual strategies of the researchers can't ensure secure storage over a longer time period and also don't support publication of the data as requested by the funding agencies.

The repository should support the mentioned data formats and also enable upload of additional description and materials related to published data.

In following document, standards, which should be taken into account, will be summarized, as well as existing benchmarking projects and related initiatives and projects.



## Appendix

### Survey text

**SURVEY to “Open Access to Research Data @ IST Austria” project**



**Please, in this survey understand the term “research data” explicitly as primary data, which are the basis for published documents (e. g. journal article, conference paper, book chapter, dissertation thesis etc.)**

1/ How are research data produced in your research process? (e.g. measuring, microscopy, taking pictures, electrophysiology traces)

2/ Out of which device/system or in which way (for non-digital data) do you get the research data?

3/ What kind of research data does your research group produce? (e.g. pictures, tables, videos, picture stacks, codes)

4/ What formats of research data does your research group produce? (e.g. JPG, TIFF, CFS, LIF format for pictures etc.)

5/ Do you have special requirements for describing the research data? (e.g. biological specimens etc.)

6/ What is the approximate amount of research data in one month?

7/ What is the approximate the average data size?

8/ Where do you store (archive) the research data once they were produced?

9/ Where do you store (archive) the research data after publication?

10/ Where do you upload (publish) your research data after publication?

11/ Do you upload (publish) all<sup>7</sup> of the research data or only a part of them after publication?

---

<sup>7</sup> All primary data, which are the basis for published documents.

## **Full list of procedures**

automated optical density measurement  
behaviour observation  
biophysical techniques  
bulk fluorescence  
collecting biological specimens from the wild  
collecting observational data of pollinators in the wild  
computer simulation  
computing results with software  
data rendering  
DNA sequencing  
downloads from public sources on the Internet  
electrophysiology  
experiments on plants  
FACS sorting/counting  
gaschromatography-mass spectrometry  
generation of simulated data  
immunological measures  
implementing algorithms  
in silico simulations, genomics  
measuring  
microscopy  
photography  
processing (feature extraction) using self-written software tools  
quantitation  
quantitative measurements  
reuse of collaborators data  
running simulations  
single cell microscopy  
software tools  
source code collecting from open source software  
spatial data  
taking pictures  
writing code

**Full list of devices/systems**

biophysical measurement tool  
CCD camera  
cluster  
datalogger  
field  
fluorescence reader  
gas-chromatograph  
GPS  
microscope  
online genomic databases  
patch clamp amplifier  
PC  
photo camera  
plate reader  
portable hard drive  
public sources from internet  
Real Time PCR machine  
spinning disc  
video camera

**List of data formats**

<b>ab1</b>	ABIF file format / highly specialized DNA sequencing format
<b>avi</b>	Audio Video Interleave movie file
<b>bmp</b>	standard Windows bitmap image
<b>c</b>	C/C++ main source code file format
<b>cfs</b>	binary format
<b>clc</b>	UltimaCalc text file
<b>cpp</b>	C++ main source code file format
<b>csv</b>	Comma Separated Value file
<b>eps</b>	Encapsulated PostScript file
<b>jpg</b>	JPEG bitmap image format file
<b>lif</b>	Leica image file
<b>ml</b>	ML language source code file
<b>mov</b>	Apple QuickTime digital movie file
<b>mp4</b>	MPEG-4 video file format
<b>mpg</b>	video file format
<b>nb</b>	Mathematica notebook file
<b>obj</b>	3D object file
<b>pdf</b>	Adobe Portable document format
<b>ps</b>	Adobe PostScript file
<b>py</b>	Python script language source code
<b>r</b>	source code
<b>sav</b>	SPSS Data Sets database file
<b>spv</b>	SPSS (Statistical Package for the Social Sciences) analysis output file
<b>stk</b>	HyperStudio stacks file
<b>tex</b>	TeX/LaTeX text document
<b>tex</b>	TeX/LaTeX text document
<b>tiff</b>	Tagged Image File Format bitmap image
<b>tikz</b>	PGF/TikZ is a tandem of languages for producing vector graphics from a geometric/algebraic description
<b>txt</b>	simple text file
<b>xlsx</b>	Microsoft Excel 2007/2010 Open XML workbook file
<b>xml</b>	XML document file