# Requirements for Persistent Identifiers within the MPI-CBG

Jeff Oegema
Max Planck Institute of Molecular Cell Biology and Genetics
eScience Seminar, March 28th, 2008

# The Question

PIDs are currently used:

- for publications (ex: CrossRef)
- within archives and libraries
- as a tool for primary scientific data (World Data Center for Climate)
- in a variety of other ways

What benefits can PIDs bring to the types of data that the MPI-CBG generates?

# Outline

1) MPI-CBG Data "Objects" and Examples

2) Current State of Data Handling and Software within the MPI-CBG

3) Benefits from PIDs and potential applications and questions

Microscopy Data - Fluorescent, Confocal, High-Throughput Screening, Electron Microscopy

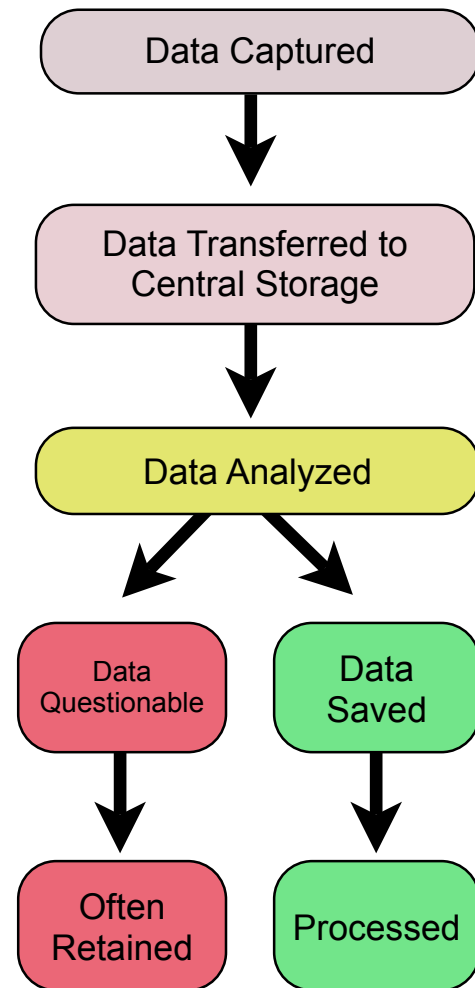Other Image Data - Protein Blots, Electrophoresis Gels, etc.

Sequence Data - DNA Sequencing Results

Mass Spec Data- Spectra and Numerical Data Files

Protocols and Methods

The MPI-CBG is heavily focused on microscopy and imaging which accounts for **95%** of our data by volume

# MPI-CBG Data Example

**Data Captured**

In most / many cases data captured is an image or movie as a result from an experiment

**Data Transferred to Central Storage**

Structure is managed by the scientist or research group - file server space allocated by project

**Data Analyzed**

Data is analyzed and evaluated, often manually

**Data Questionable**

**Data Saved**

After analysis the data is deemed to be valuable, a failure, or the results are indeterminate.

**Often Retained**

**Processed**

It is often unclear for quite some time after capture if the data has meaning or value.

# What does the data look like?
## Primary Data



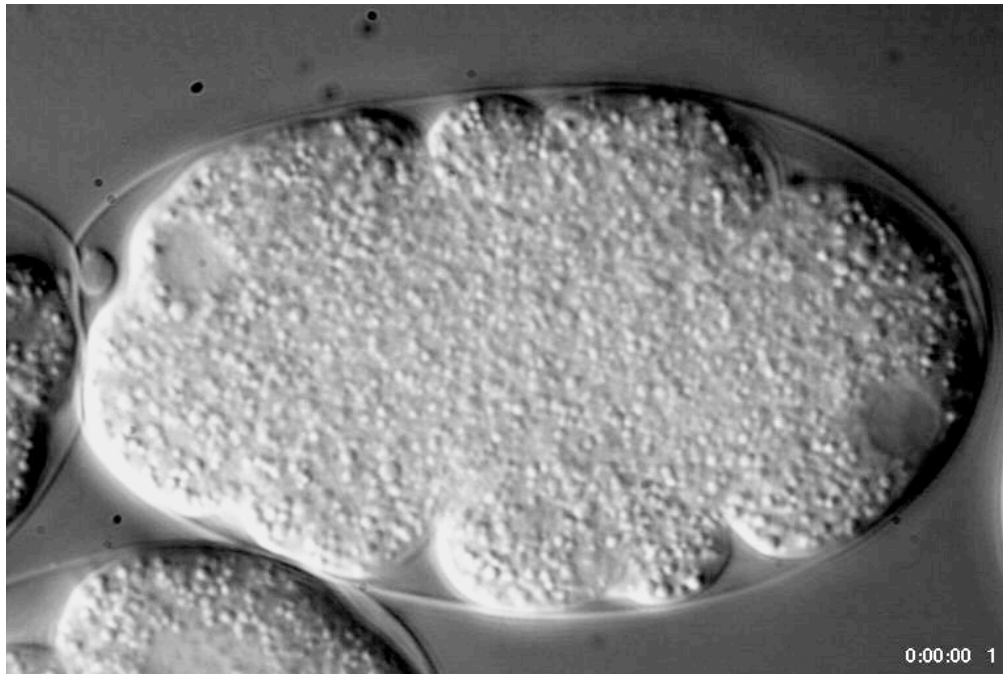In Zeiss LSM file format (Zeiss Meta Imaging System)

Readable using the Zeiss software or an ImageJ Plugin

This image is of larval tissue of a fruit fly and shows the expression of the Dally protein in green and lipoprotein in red.

Dally is overexpressed.

Christina Eugster - Eaton Lab

# What does the data look like?
## Primary Data



C. Elegans Division - Hyman Lab

Captured using normal contrast microscopy

Uncompressed movie format

No metadata encoded with file on capture, but was manually added

The movie shows the merging of the egg and sperm cells and the first divisions of a C. Elegans embryo

# What does the data look like?
## Primary Data



0:00:00  1
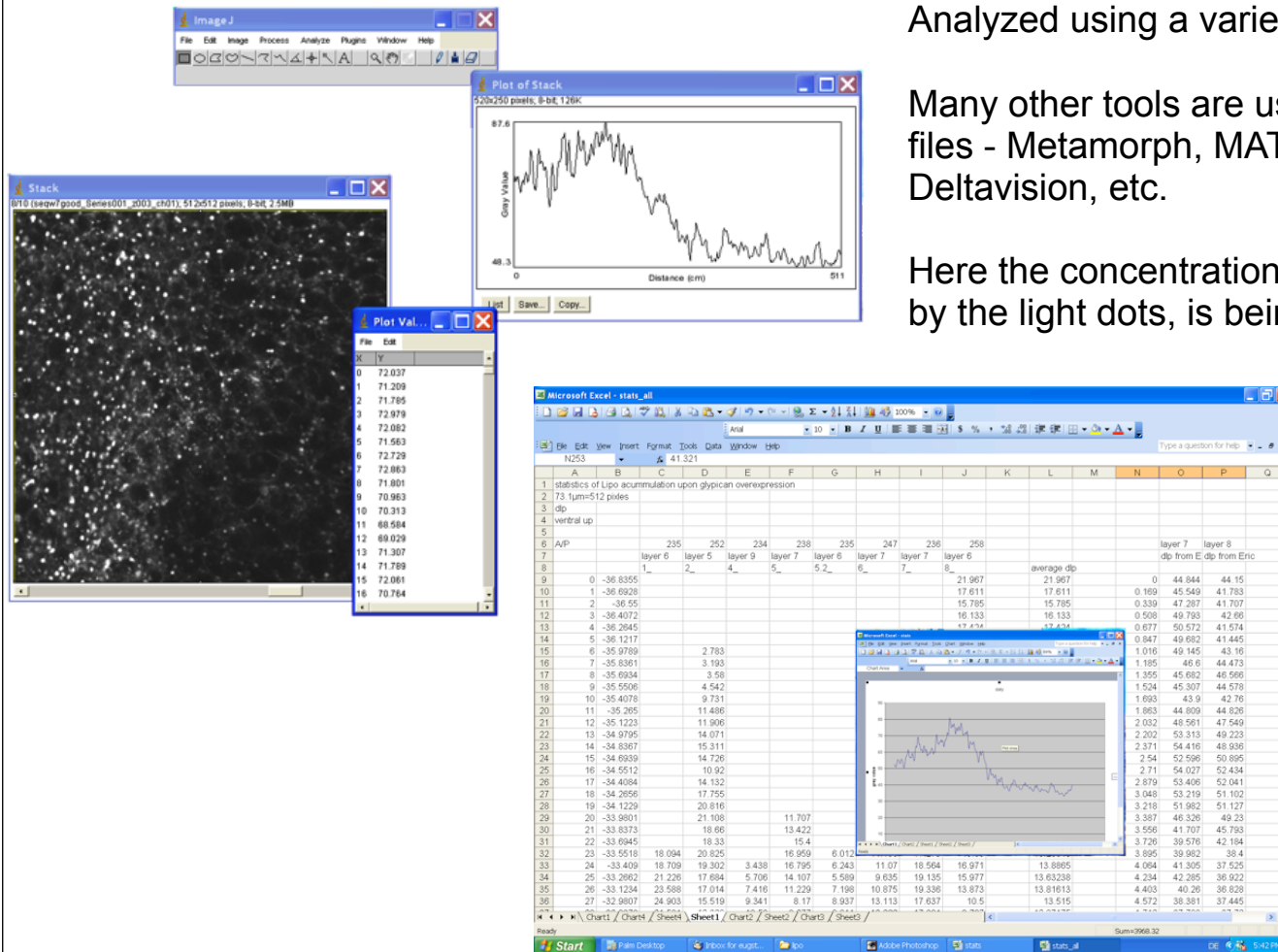
C. Elegans Division - Hyman Lab

This movie shows the failure pattern when gene H04J21.3 was knocked down.

This pattern was interpreted by eye and was classified as a "Spindle Assembly" problem

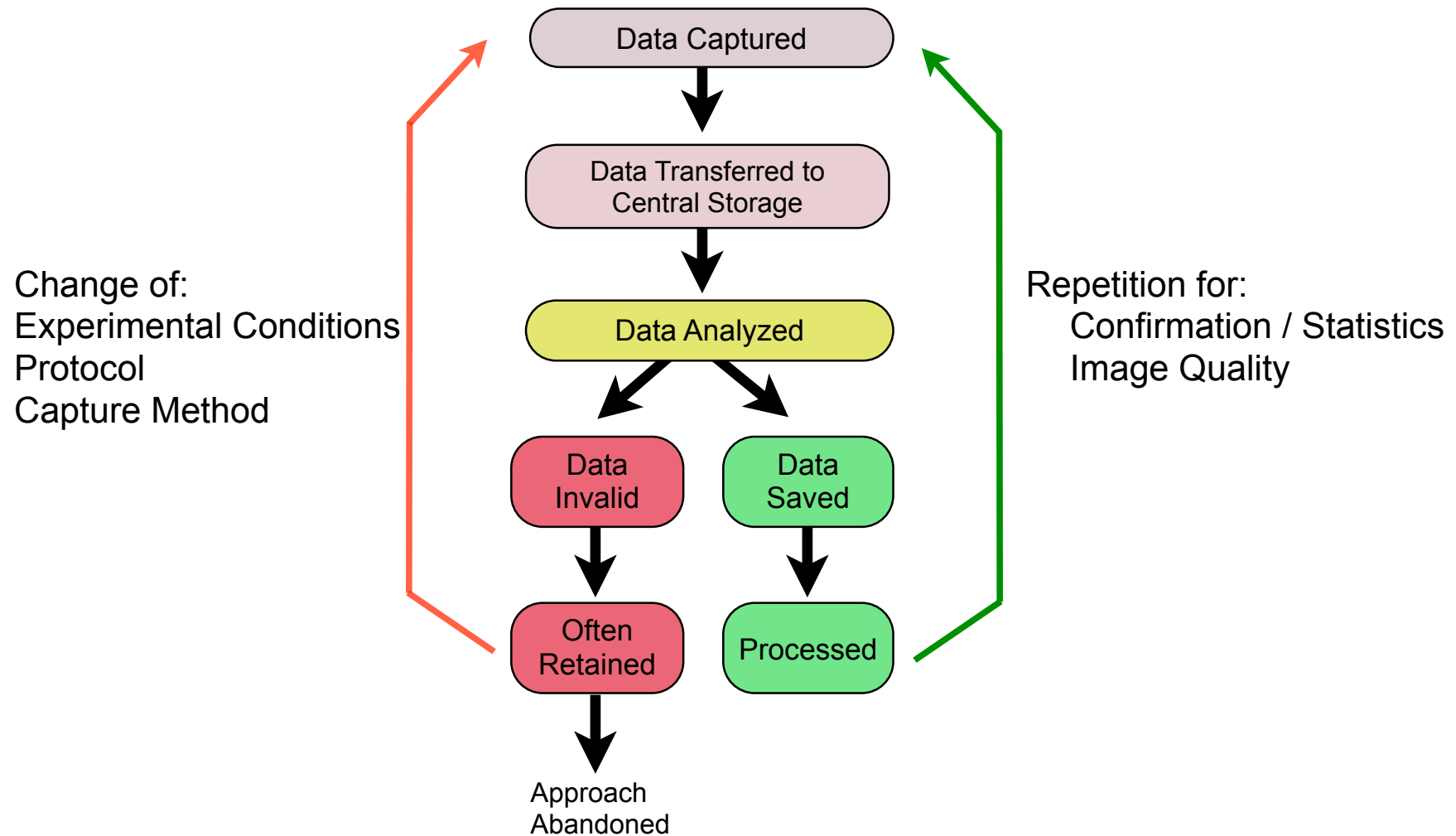# What does the data look like?
## Analysis of Data



Analyzed using a variety of tools - Example ImageJ

Many other tools are used and generate a variety of files - Metamorph, MATLAB, Huygens, Definiens, Deltavision, etc.
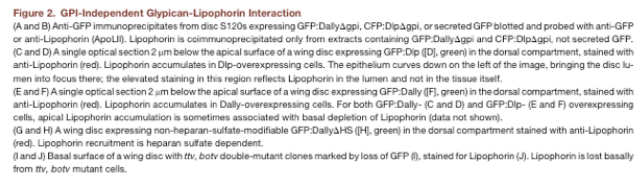
Here the concentration of lipoprotein, as indicated by the light dots, is being examined and plotted.

# What does the data look like?
## Figures and Publication



Multiple experiments and primary data sets come together for a single publication figure.

Here the concentrations of several proteins are being shown side by side.

**Developmental Cell**
Lipoprotein-Heparan Sulfate Interactions

Figure 2. GPI-Independent Glypican-Lipophorin Interaction
(A and B) Anti-GFP immunoprecipitates from disc S120s expressing GFP:DallyΔgpi, CFP:DlpΔgpi, or secreted GFP blotted and probed with anti-GFP or anti-Lipophorin (ApoLII). Lipophorin is coimmunoprecipitated only from extracts containing GFP:DallyΔgpi and CFP:DlpΔgpi, not secreted GFP.
(C and D) A single optical section 2 μm below the apical surface of a wing disc expressing GFP:Dlp (D), green) in the dorsal compartment, stained with anti-Lipophorin (red). Lipophorin accumulates in Dlp-overexpressing cells. The epithelium curves down on the left of the image, bringing the disc lumen into focus there; the elevated staining in this region reflects Lipophorin in the lumen and not in the tissue itself.
(E and F) A single optical section 2 μm below the apical surface of a wing disc expressing GFP:Dally ([F], green) in the dorsal compartment, stained with anti-Lipophorin (red). Lipophorin accumulates in Dally-overexpressing cells. For both GFP:Dally- (C and D) and GFP:Dlp- (E and F) overexpressing cells, apical Lipophorin accumulation is sometimes associated with basal depletion of Lipophorin (data not shown).
(G and H) A wing disc expressing non-heparan-sulfate-modifiable GFP:DallyΔHS ([H], green) in the dorsal compartment stained with anti-Lipophorin (red). Lipophorin recruitment is heparan sulfate dependent.
(I and J) Basal surface of a wing disc with ttv, botv double-mutant clones marked by loss of GFP (I), stained for Lipophorin (J). Lipophorin is lost basally from ttv, botv mutant cells.

# Complizations



Protocols change, sometimes with every repetition

A proper (and complete) description must be bound to each data set or the data is uninterpretable

Automation of Data Generation:

• Simplifies protocol data as one protocol is used for each sample

• Produces millions of image and analysis results



• Data collection is automated

• Analysis may be automated

Image courtesy of Tecan Group Ltd.

**Data Capture**

19.1 TB of Image Data (tiff), over 2,000,000 images

**Analysis**

Thousands of computer hours were used to generate numerical analysis

**Interpretation**

Much of the data still human interpreted based on the numerical analysis

**Analysis is only as good as the image processing tools used**

**Human interpretation is based on current knowledge and limited by human perception.**

# Data Summary

## Experimental Data

Generated by a specific protocol

Meant to answer a single question

May not have additional information

Generally is image or movie data

Analysis generates numerical data

Almost always repeated, sometimes with protocol changes



Lipoprotein Distribution - Eaton Lab

## Screening Data

Before going to screening the protocol must be locked

Experimentation is automated

Analysis is computationally intensive

Still usually requires human interpretation

Is likely to be useful for future reanalysis



GWS Imaging - Zerial Lab

## Protocols and Methods

A protocol is linked to an experiment, or to a screen

The format of the protocols kept varies highly

Experimental data is only meaningful in combination with the protocol

Protocols themselves have value as they teach methods to scientists



Staining Protocol - Eaton Lab

# Data Summary

Data within the MPI-CBG is captured and organized in many ways.

Before we can present primary data to the world we must improve our own internal data handling.

Data within the building must be consistent, uniquely identified, and associated with the metadata.

If PIDs are ready for us, are we ready for them?

Improvement of internal standards

Professionalization of software development for data management (including internal metadata and identifiers)

Data management migrating out of individual non-standard efforts towards standards-compliant systems provided by a central software development group

Formation of an image processing and analysis facility to provide central know-how and tools

# Current Efforts

Specific "display" database systems provide access to information for a given screen or publication when this is considered necessary

Data management systems are being created to assist with information flow and experiment management

Protocols are currently managed by wiki, in a simple database system, or manually, depending on group. Plans exist to standardize format, structure, and location of these protocols.

# Genome Wide Screen Display Database

**CBG**
Max Planck Institute
of Molecular Cell Biology
and Genetics

Administrator Menu | **User Menu**

Gene Scores

**15E1.2**                                                    Oligo Details

**Gene Profile**

| | |
|---|---|
| Numb. Ves. (Channel 1) Mask = TRUE | 4.13793 |
| Total Intens. (Channel 1) Mask = TRUE | -0.67582 |
| Integ. Ves. Intens. (Channel 1) Mask = TRUE | 0.383962 |
| Mean Area (Channel 1) Weighed=TRUE WeighingFunc=GetVolume CalcType = Mean | -2.73381 |
| Mean Area (Channel 1) Weighed=TRUE WeighingFunc=GetMeanIntensity CalcType = Mean | -2.84289 |
| Mean Area (Channel 1) Weighed=FALSE CalcType = Median | -3.09261 |
| Mean Elongation (Channel 1) Weighed=TRUE WeighingFunc=GetVolume CalcType = Mean | -0.578184 |
| Mean Elongation (Channel 1) Weighed=TRUE WeighingFunc=GetMeanIntensity CalcType = Mean | -1.36949 |
| Total Intens. (Channel 1) Mask = TRUE | -0.67582 |
| Integ. Ves. Intens. (Channel 1) Mask = TRUE | 0.383962 |
| Mean Area (Channel 1) Weighed=TRUE WeighingFunc=GetVolume CalcType = Mean | -2.73381 |
| Mean Area (Channel 1) Weighed=TRUE WeighingFunc=GetMeanIntensity CalcType = Mean | -2.84289 |
| Mean Area (Channel 1) Weighed=FALSE CalcType = Median | -3.09261 |

# CBG

Max Planck Institute
of Molecular Cell Biology
and Genetics

## Huttner Laboratory Storage Database

| Administrator | Main | Search | Import |

Antibody    DNA Construct    SiRNA    GMO    Oligo    Protocol    Chemical    Consumables    Supplier

## List Protocol

17 items found, displaying 1 to 15.[First/Prev] **1**, 2 [Next/Last]

| ⬦ Name | ⬦ File | ⬦ Comment | ⬦ Bacteria Transformation | ⬦ Submitted By |
|---|---|---|---|---|
| BrdU staining | BrdU – Labelling.doc | | false | haffner |
| Collagen volumes table | Collagen_mix_volumes.xls | To prepare smaller amounts of 1.5 mg/ml collagen for slice culture. | false | mora |
| esiRNA preparation | esiRNA.doc | | false | marzesco |
| grids for microinjection (Warner Instruments) | | 0.5x0.5cm and 1x1cm grids, nylon threads. These grids are usualy suitable to hold slices during electrophisiology recordings. | false | taverna |
| grids for microinjection (workshop) | | 1cm x1cm grids, nylon threads These grids are used to hold slices during microinjection. These grids can accomodate 250-300micron-thick slices. | false | taverna |
| IF staining (Yoichi's protocol) | Immunostaining (YK).doc | | false | haffner |
| Immunofluorescence protocol_AMM | immunofluorescence.doc | | false | marzesco |
| In utero electroporation | in utero EP.pdf | from Tetsuichiro Saito web page. | false | marzesco |
| Mowiol | MOWIOL.doc | | false | haffner |
| Nissl staining (cresyl violet) | Cresyl_violet_protocol.doc | Nissl staining for sections on glass slides | false | pulvers |
| Paraffin embedding and deparaffinization | Protocol for paraffin embedding and deparaffinization.doc | | false | fietz |
| Sequencing facility primer list | Primer_List.pdf | Sequencing primers provided by sequencing facility | false | pulvers |
| Slice Culture Medium | Slice Culture Medium.doc | | false | taverna |
| staining with boiling | staining with boiling.doc | e.g.Pax6,Tbr1,Tbr2 | false | haffner |
| staining with boiling and HCL-treatment | staining with boiling and HCL.doc | e.g.Pax6,Tbr2,Tbr1,BrdU,GFP | false | haffner |

Export options: CSV I Excel I XML I PDF

| Add Protocol |

# DNA Sequencing Request Database

**CBG**
Max Planck Institute
of Molecular Cell Biology
and Genetics

| Administrator Menu | **Technician Menu** | User Menu |

Kits    Primers    Templates    **All Sequencing Plates**    Uncompleted Sequencing Plates    Unassigned Order Entries

**Edit Plate**

|     | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| --- | - | - | - | - | - | - | - | - | - | -- | -- | -- |
| A | ■ | ■ | ■ | ■ | □ | □ | □ | □ | □ | □ | □ | □ |
| B | ■ | ■ | ■ | ■ | □ | □ | □ | □ | □ | □ | □ | □ |
| C | ■ | ■ | ■ | □ | □ | □ | □ | □ | □ | □ | □ | □ |
| D | ■ | ■ | ■ | □ | □ | □ | □ | □ | □ | □ | □ | ■ |
| E | ■ | ■ | □ | □ | □ | □ | □ | □ | □ | □ | □ | ■ |
| F | ■ | ■ | □ | □ | □ | □ | □ | □ | □ | □ | □ | ■ |
| G | ■ | ■ | □ | □ | □ | □ | □ | □ | □ | □ | □ | ■ |
| H | ■ | ■ | ■ | □ | □ | □ | □ | □ | □ | □ | □ | ■ |

**Plate** `seq_080110_A`

| Redo | Well | UEN | Name | Customer | Billing Group | Date Submitted | Date Completed | Template variance | Primer | Primer Tm | Kit | Read Length | Template | Tann | Comment |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| ☐ | A01 | 1 | 1 | Seidel | TILLING | 10.01.2008 15.35 | | ☐ | M13 reverse | 56 | 1/5 | | PCR | 52 | |
| ☐ | B01 | 2 | 1a | Seidel | TILLING | 10.01.2008 15.35 | | ☐ | M13 universal (-21) | 54 | 1/5 | | PCR | 52 | |
| ☐ | C01 | 3 | 1b | Seidel | TILLING | 10.01.2008 15.35 | | ☐ | T3 | 54 | 1/5 | | PCR | 52 | |
| ☐ | D01 | 4 | 1c | Seidel | TILLING | 10.01.2008 15.35 | | ☐ | T7 promoter | 56 | 1/5 | | PCR | 52 | |
| ☐ | E01 | 5 | 2 | Seidel | TILLING | 10.01.2008 15.35 | | ☐ | M13 reverse | 56 | 1/5 | | PCR | 52 | |
| ☐ | F01 | 6 | 2a | Seidel | TILLING | 10.01.2008 15.35 | | ☐ | M13 universal (-21) | 54 | 1/5 | | PCR | 52 | |
| ☐ | G01 | 7 | 2b | Seidel | TILLING | 10.01.2008 15.35 | | ☐ | T3 | 54 | 1/5 | | PCR | 52 | |
| ☐ | H01 | 8 | 2c | Seidel | TILLING | 10.01.2008 15.35 | | ☐ | T7 promoter | 56 | 1/5 | | PCR | 52 | |

# PID Questions

The questions then become-

What role can and should PIDs play in our environment as we move towards internal standardization?

What role should they play in our distribution of data to a wider audience?

# PID Questions

- **Granularity** - At what level should a dataset be given a PID?

- **Timing** - When in the data life cycle should a dataset be given a PID?  What do you do with multiple versions?

- **Access** - Should data at preliminary stages be accessible to the world or even to other internal groups?  What should be done to address data privacy concerns, and data privacy requirements that vary over time?

- **Data Lockdown** - After a PID is assigned, how much can the data be changed?

- **Metadata** - What information should be associated with our PIDs, and are these schemas already existing or do we have to invent them?

# Data Search

Search | All | Lipoprotein

- All
- Analyzed Data
- Primary Data
- Protocols
- Publications

Search    Reset

## Publications (2)

| Title | Authors | Date of Pub | |
|---|---|---|---|
| Lipoprotein-heparan sulfate interactions in the Hh pathway. | Eugster C, Panáková D, Mahmoud A, Eaton S. | July, 2007 | View Detail |
| Lipoprotein particles are required for Hedgehog and Wingless signalling | Panakova, Daniela; Sprong, Hein; Marois, Eric; Thiele, Christoph; Eaton, Suzanne | June, 2005 | View Detail |

## Analyzed Data Sets (5))

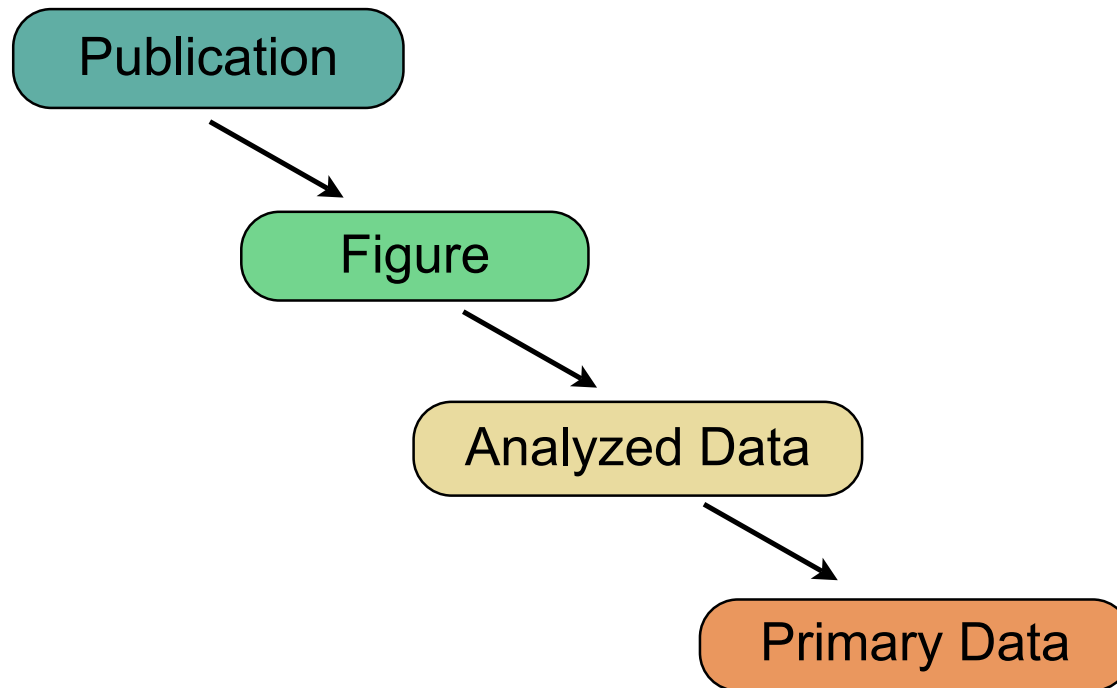| Title | Description | Format | Date Analyzed | |
|---|---|---|---|---|
| QI_25-5-2007-Anti_Lipo | Quantification of fluorescence intensity of Lipophorin | XLS File | 15-6-2007 | View Detail |
| CL-27-7-2007-Anti_Lipo_Hh_Ptc | Co-localization of Lipophorin, Hedgehog, and Patched with Image J | XLS FIle | 12-8-2007 | View Detail |

## Primary Data Sets (15)

| Title | Description | Format | Date Captured | |
|---|---|---|---|---|
| 25-5-2007-Anti_Lipo | Overexpression of Dally, Dally Like, Dally Delta HS stained for Anti Lipophorin | LSM Database | 25-5-2007 | View Detail |
| 27-7-2007-Anti_Lipo_Hh_Ptc | Overexpression of Dally, Dally Delta GPI, stained for Anti Lipophorin, Hedgehog, and Patched | LSM Database | 27-7-2007 | View Detail |

## Protocols (3)

| Title | Description | Format | Date Finalized | |
|---|---|---|---|---|
| Co-immunoprecipitation of Lipoprotein Particles | Protocol for the Co-immunoprecipitation of Lipoprotein Particles | PDF | 20-11-2006 | View Detail |
| Fluorescent Labeling of Lipoprotein Particles | Protocol for the Fluorescent Labelling of Lipoprotein Particles | PDF | 15-6-2006 | View Detail |

# Publication Drill-Down



A scientist reading a publication has access to the primary data and can verify the correctness of the publication

# Archiving and Data Loss

Scientist arrives at the MPI-CBG
(Student, Post-Doc)

↓

Scientist gets project space
uses it to store data

↓

Scientist leaves, transfers some
data and knowledge to the group

↓

Project space is closed
Data is transferred to tape

↓

All significant data is stored but
is practically "unrecoverable"

# Conclusion

**Benefits**

Screening Data has the potential for future data mining to yield significant gain

Linking analyzed and primary data to publications allows greater transparency and validation of scientific work

Successful metadata binding to experimental data prevents data loss, allows greater reuse, and easier automated data mining

PIDs allow easier reuse of data at a later point, greater compatibility/interoperability with other systems and potentially other organizations with related data.

**Difficulties**

It is difficult to determine at what granularity and when to assign PIDs.  This for us will require significant work on data organization structure.

We have only light experience with PIDs as a concept in-house at the moment.

Any system requiring the addition of metadata to primary research results must be extraordinarily easy to use or the research scientist won't use it

Resources and funding within the MPI-CBG for such projects is difficult to obtain as the software facility and development resources are stretched with existing load.

# Conclusion

The MPI-CBG is currently solidifying its internal data structures.

Until these structures are solid, its difficult to expand to greater world release of digital data through PIDs

The concept of PIDs and the associated metadata seems to have potential benefits for they types of data the MPI-CBG handles

The MPI-CBG will revisit the concept of PIDs as it's software development progresses and appreciate feedback from the community.

Its likely there are tools and resources available that we are unaware of

# Thanks!

**Software Engineering Facility**
Tim Cross - cross@mpi-cbg.de

**Image Processing and Analysis Facility**
Jean-Yves Tinevez - tinevez@mpi-cbg.de

**IT Coordination**
Jeffrey Oegema - joegema@mpi-cbg.de