

# Image Data Storage

Josh Moore  
University of Dundee

The OME Consortium  
*openmicroscopy.org*  
*@openmicroscopy*



University  
of Dundee



# Outline

- OME definition of Image

- File-formats

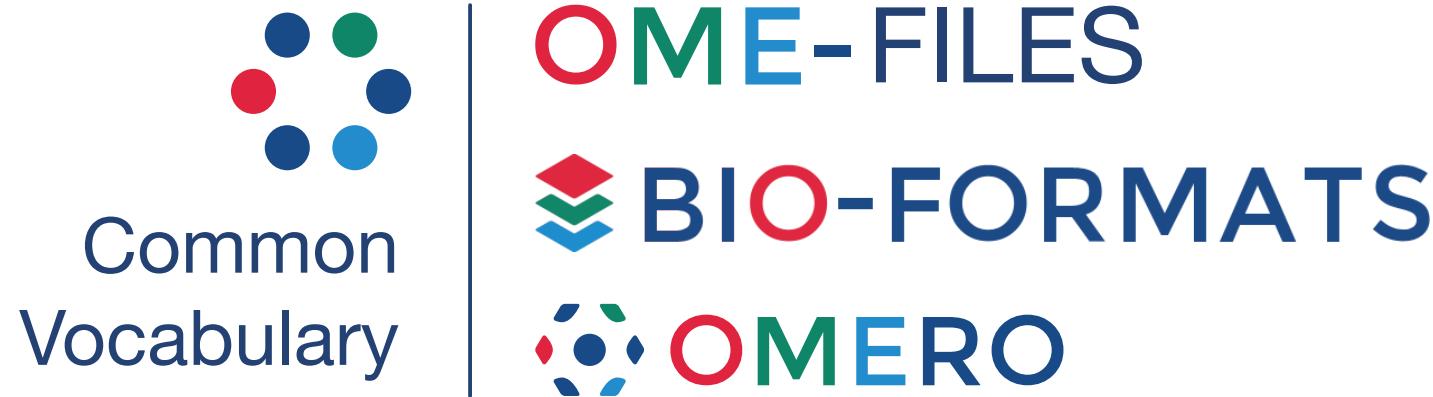
- Current storage
  - Challenges, etc.

- Object storage

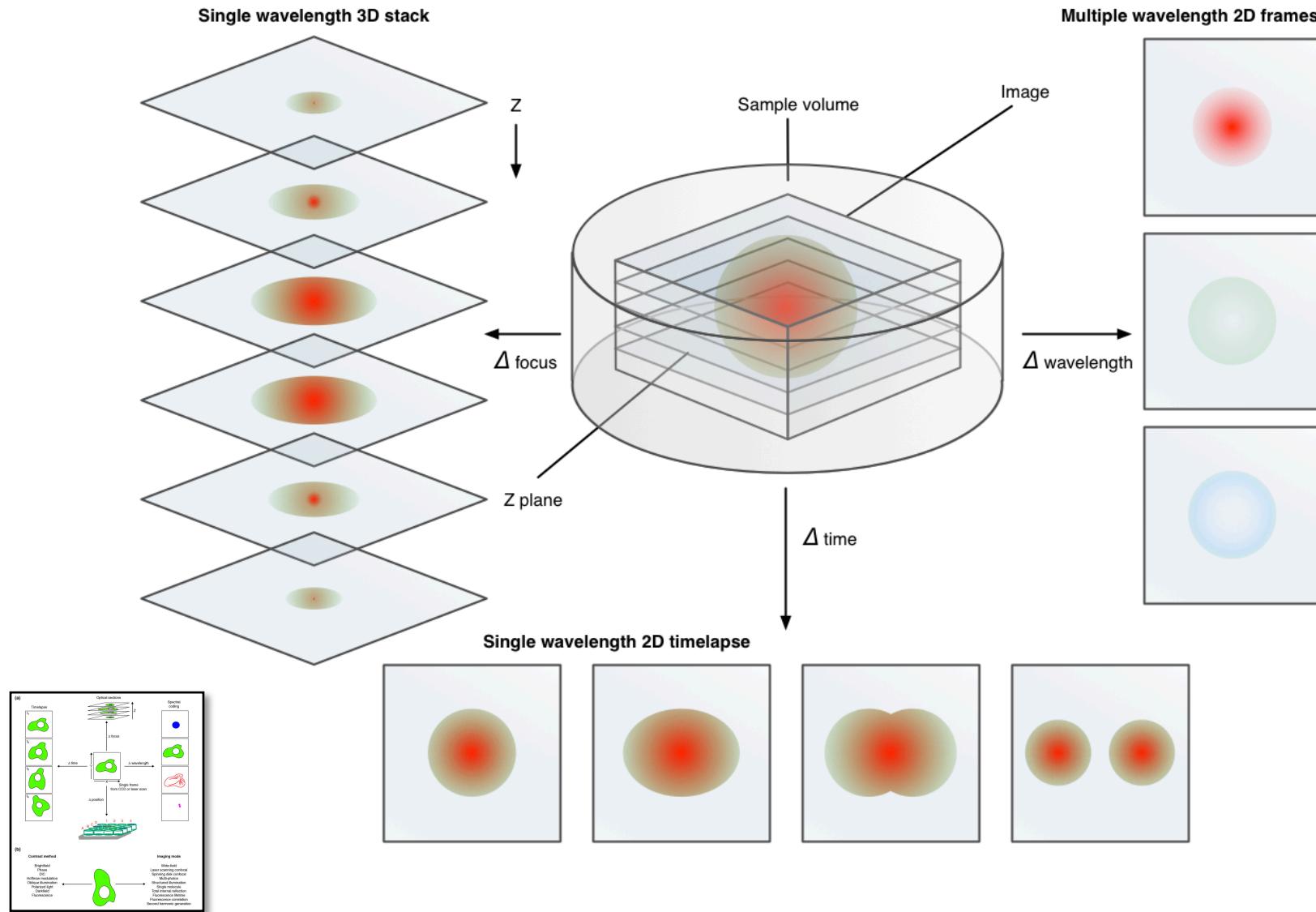
- Expectations
  - Cost/Value
  - Requirements

# IMAGES

# OME Strategy

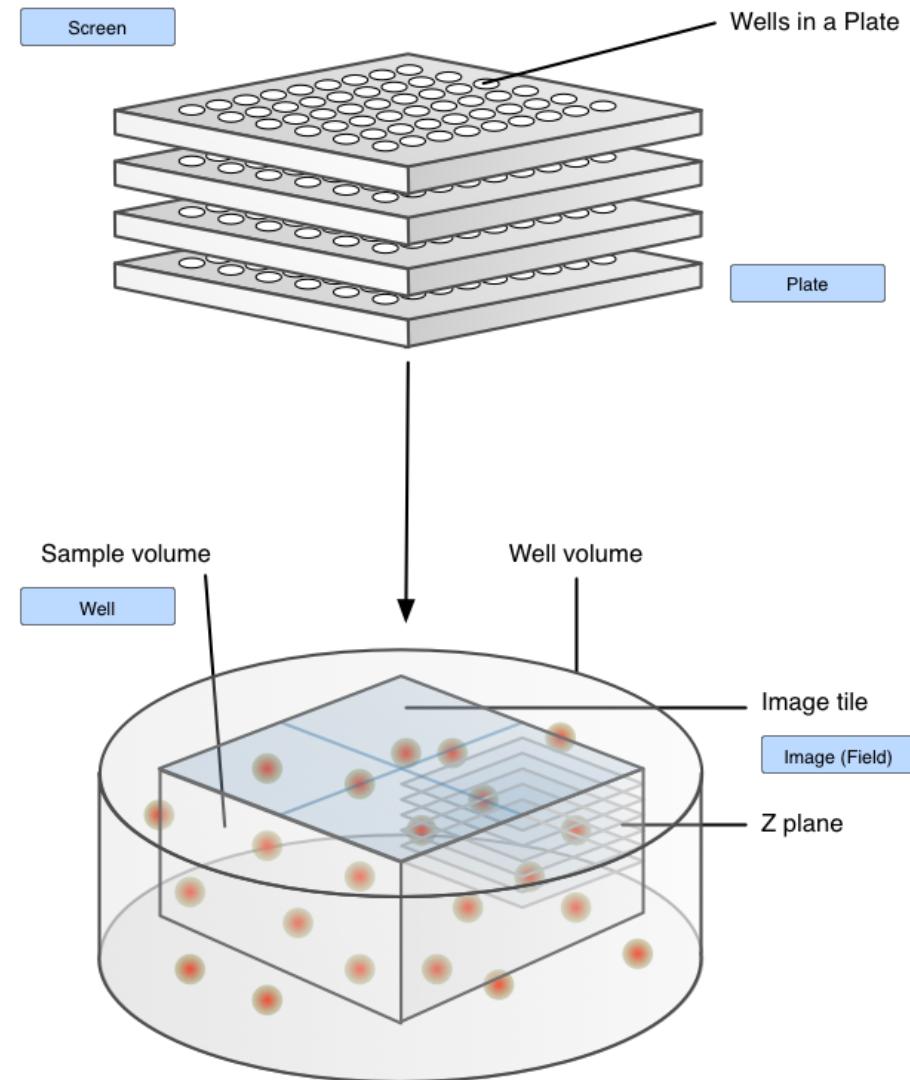


# Spatio-temporal dimensions + measurements



<https://genomebiology.biomedcentral.com/articles/10.1186/gb-2005-6-5-r47>

# High-content screening dimensions



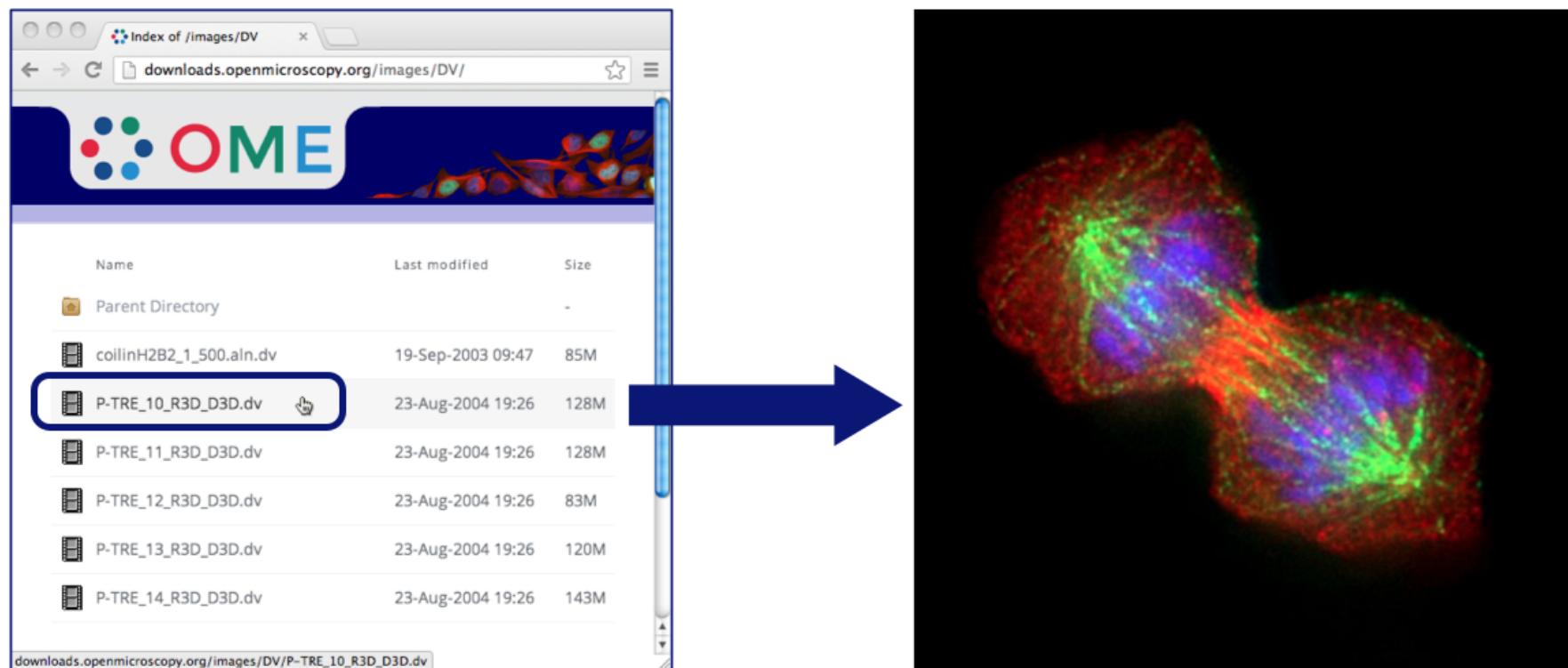
# FILE-FORMATS

# Bio-Formats

140+ proprietary formats

Format	Extensions	Pixels	Metadata	Openness	Presence	Utility	Export	BSD	Multiple Images	Pyramid
3i SlideBook	.sld	▲	▼	▼	▼	▼	✗	✗	✓	✗
Andor Pro-Imaging Revision (ABD) TIFF	.tif	▲	▲	▼	▼	▼	✗	✗	✗	✗
AIM	.aim	▲	▼	▲	▲	▼	✗	✗	✗	✗
Alicona SD	.aisd	▲	▲	▲	▲	▲	✗	✗	✗	✗
Varian FDF	.fdf	■	▼	▼	▼	▼	✗	✗	✗	✗
Vecco AFM	.hdf	■	▼	▲	▼	■	✗	✗	✗	✗
VG SAM	.dti	■	▼	▼	▼	▼	✗	✗	✗	✗
VisiTech XYS	.xys, .html	▲	■	▼	▼	■	✗	✗	✓	✗
Velocity	.mvd2	■	■	▼	▼	▼	✗	✗	✓	✗
Velocity Library Clipping	.acff	■	■	▼	▼	▼	✗	✗	✗	✗
WA-TOP	.wat	■	▼	▼	▼	▼	✗	✗	✗	✗
Windows Bitmap	.bmp	▲	▲	▼	▼	▼	✗	✓	✗	✗
Woolz	.wlz	▲	▼	▲	▼	▼	✓	✗	✗	✗
Zeiss Axio CSM	.lms	■	▼	▼	▼	▼	✗	✗	✗	✗
Zeiss AxioVision TIFF	.xml, .tiff	▲	▲	■	▼	▼	✗	✗	✓	✗
Zeiss AxioVision ZVI (Zeiss Vision Image)	.zvi	▲	▲	▲	■	■	✗	✗	✗	✗
Zeiss CZI	.czi	▲	▲	▲	▼	■	✗	✗	✓	✓
Zeiss LSM (Laser Scanning Microscope) 510/710	.lsm, .mdb	▲	▲	■	■	■	✗	✗	✓	✗

## Filesets: 1-1



## Filesets: 1-N

Index of /images/SVS

downloads.openmicroscopy.org/images/SVS/

**OME**

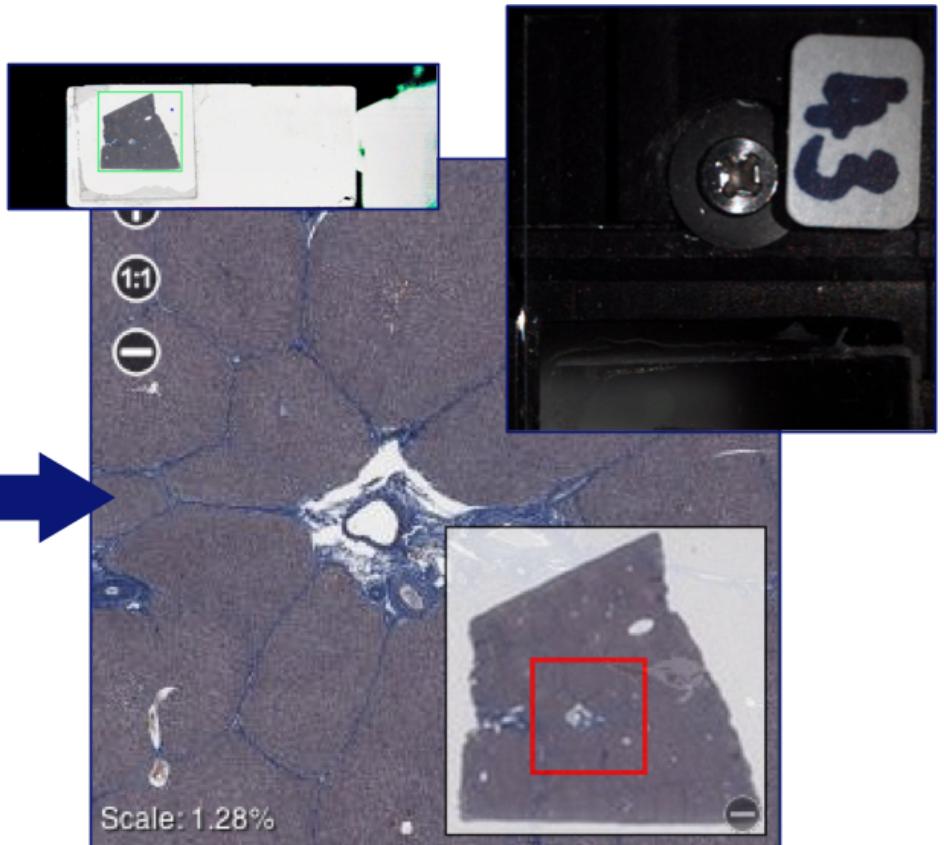
Name	Last modified	Size
Parent Directory		-
77917.svs	08-May-2013 20:48	671M
<b>77928.svs</b>	08-May-2013 20:47	581M

Listing styled using Apaxy by @adamwhitcroft

© 2000-2013 University of Dundee & Open Microscopy Environment. Creative Commons Attribution 3.0 Unported License

OME source code is available under the GNU General public license or through commercial license from Glencoe Software

<downloads.openmicroscopy.org/images/SVS/77928.svs>



# Filesets: N-N

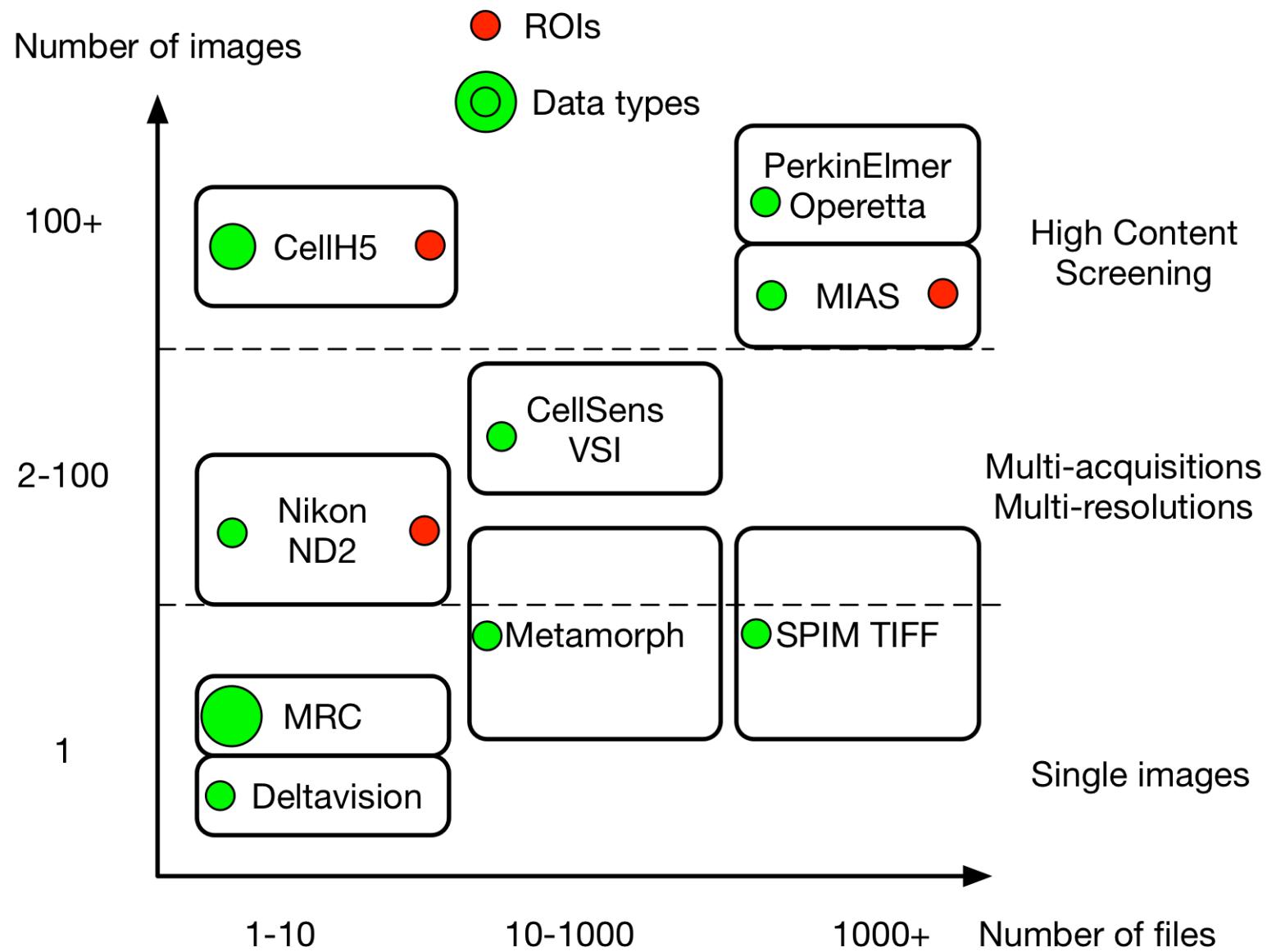
The image shows two side-by-side views of a microscopy dataset. On the left, a file browser window titled "Index of /images/HCS/Op..." displays a list of files. The columns are "Name", "Last modified", and "Size". The files listed are:

Name	Last modified	Size
Parent Directory		-
r01c01-0565849973.tif.gz	03-Apr-2013 09:17	2.0M
r01c01-0816390759.tif.gz	03-Apr-2013 09:17	1.9M
r01c01-0991679013.tif.gz	03-Apr-2013 09:17	1.7M
r01c02-0203004369.tif.gz	03-Apr-2013 09:17	2.0M
r01c02-1384679954.tif.gz	03-Apr-2013 09:17	2.0M
r01c02-1749148666.tif.gz	03-Apr-2013 09:17	1.7M
r01c03-0847230653.tif.gz	03-Apr-2013 09:19	2.1M
r01c03-1032313189.tif.gz	03-Apr-2013 09:19	1.7M
r01c03-1985691398.tif.gz	03-Apr-2013 09:19	2.2M
r01c04-0512157530.tif.gz	03-Apr-2013 09:19	2.2M
r01c04-1271778531.tif.gz	03-Apr-2013 09:19	1.7M
r01c04-1630181344.tif.gz	03-Apr-2013 09:19	2.2M
r01c05-0006079694.tif.gz	03-Apr-2013 09:21	2.2M
r01c05-0175927353.tif.gz	03-Apr-2013 09:21	1.7M
r01c05-2128047606.tif.gz	03-Apr-2013 09:21	2.2M

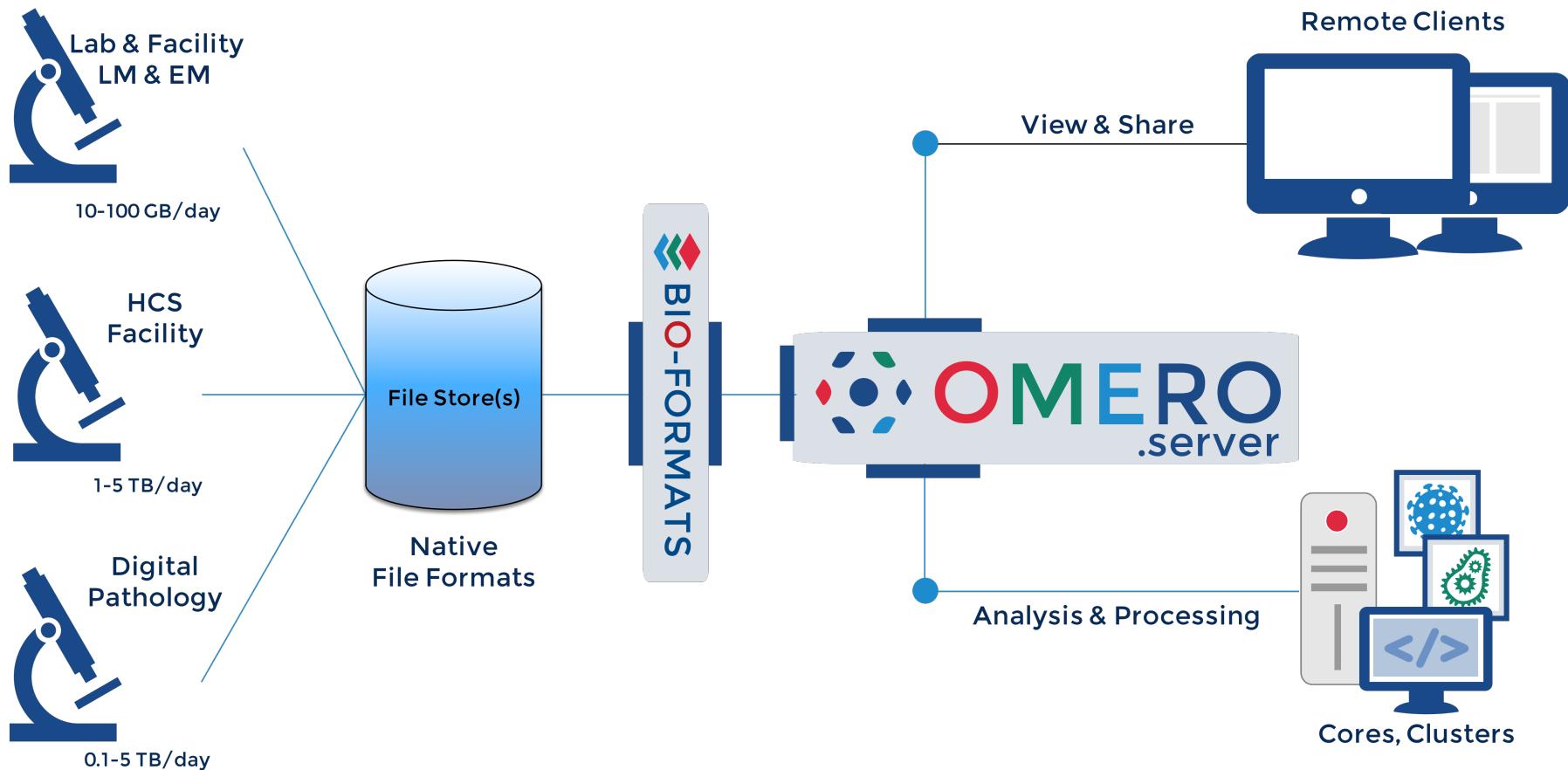
At the bottom of the list, there are three dots (...).

A large blue arrow points from the file list to the right panel, which displays a 4x6 grid of microscopy images. The grid is labeled "Field#1" at the top and has columns numbered 1 through 24. The rows are labeled with letters A through P. Each image in the grid shows a field of view with various cellular structures.

# Fileset Diversity

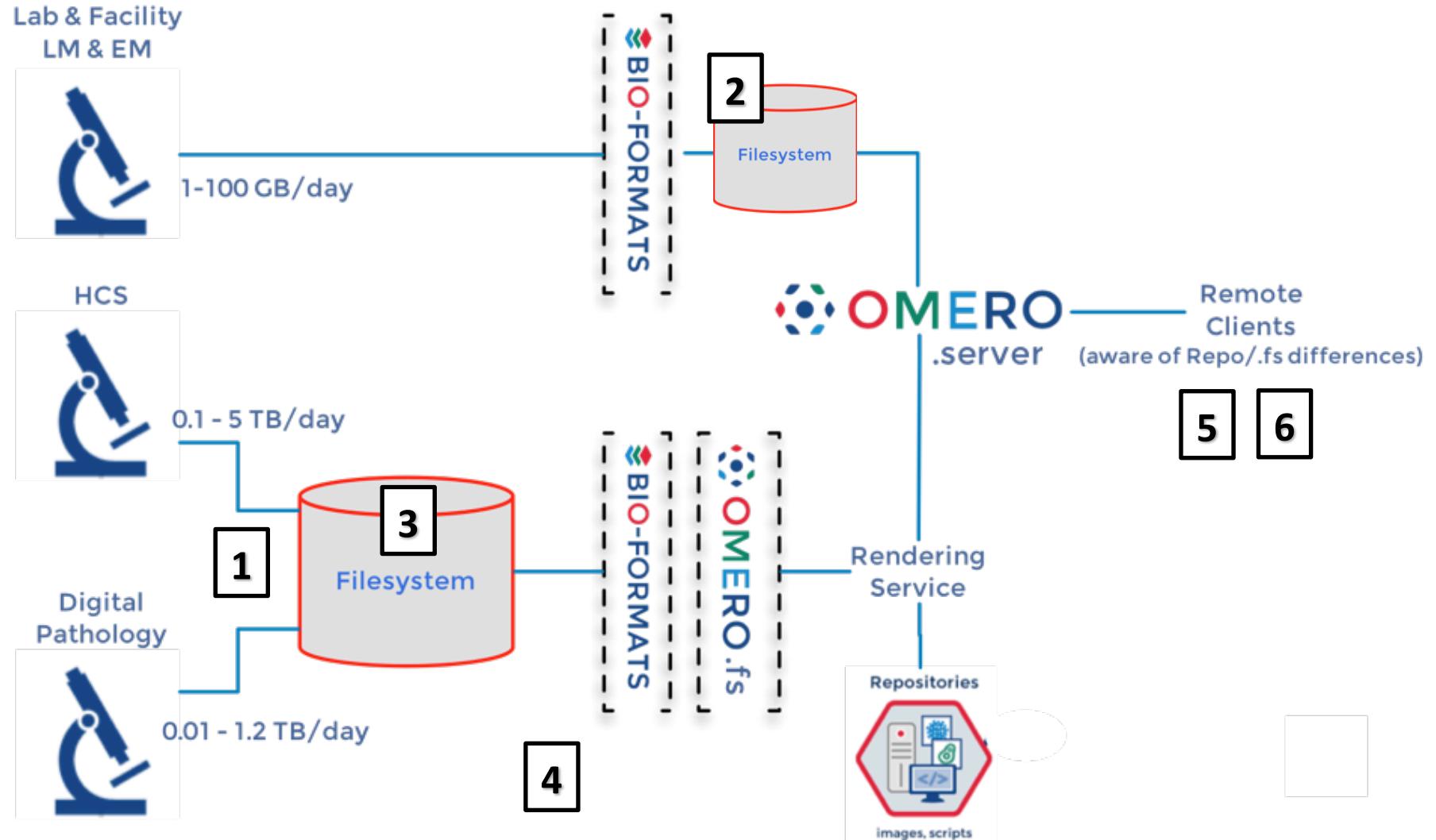


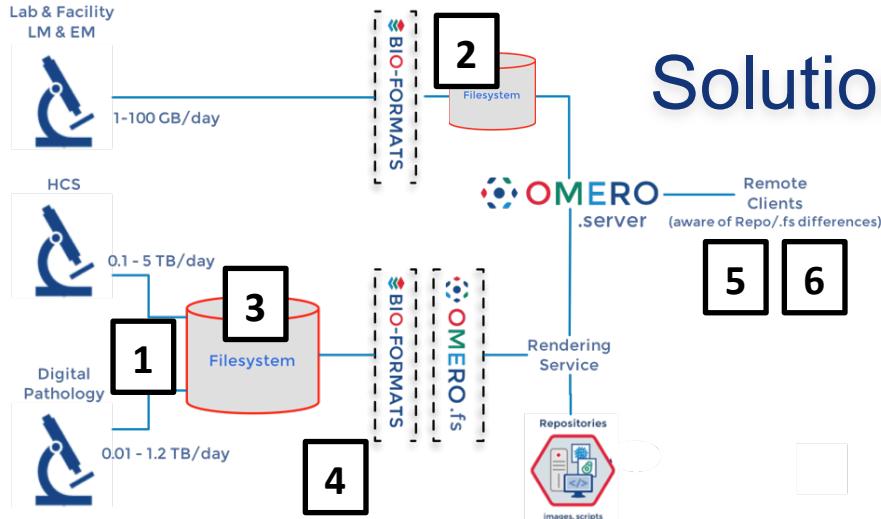
# Image Data Access: Today



# OBJECT-STORES

# Image Data Access: Past & Future





# Solution Trade-offs

		Scope	Latency	Storage	Value	ETA
1	Filesystem cache	Backend cache	Slow once	>1x	Fast & simple	Straight-forward once defined
2	Direct storage (“ROMIO”)	Swappable backend	Slow <b>always</b>	1	Space-efficient	Requires new import strategy
3	Reader/Writer	Supported format	Slow <b>always</b>	1	Stability	<b>Longest; requires spec</b>
4	Plane/chunk cache	Flexible layer	Slow once	>1x	Can be combined	Relatively simple
5	Client loads from S3	Full S3 integration	<i>Potentially fast</i>	1	Most “modern”	Requires significant changes
6	Client caches to S3	Frontend cache	Slow once	>1x	Trans-transparent	Simple

## Questions

- Will the object store be the primary storage?
  - ▶ Assuming: yes
- Will the files be deleted?
  - ▶ Assuming: yes
- Will the storage be file-based?
  - ▶ Strongly urge: no
- Will multi-resolutions be stored/generated?
  - ▶ Assuming: yes
- Will orthogonal views be needed? 2D needed?
  - ▶ Probably, but what's the default? e.g. XYZCT
- Will multi-fileset correlation be needed?

## Last thoughts

- Caching as alternate views
- Diversity of acquisition, usage & domains
- Cost of proprietary formats
- Single API vs. maintenance burden
- Relationship to file-based solutions (HDF5)

**THANKS**

# EXTRA SLIDES

# How to Source Data



Simon Li  
@crucifixkiss



Following

Another 20TB of #BigData for the  
@openmicroscopy @BBSRC @emblebi image  
repository



RETWEET  
**1**

LIKES  
**3**



1:43 PM - 15 Sep 2015



## Fileset challenges

Largest study:  
**13 TB**  
—  
idr0013  
(Neumann)

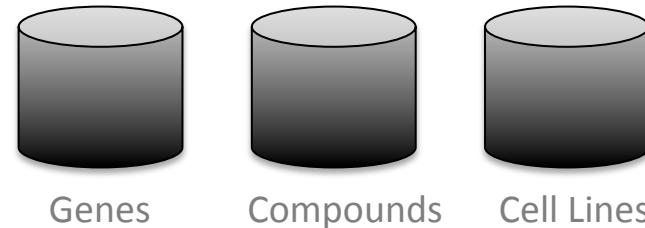
Most files:  
**5 million**  
—  
idr0016  
(Wawer)

Most wells:  
**300 K**  
—  
idr0009  
(Simpson)

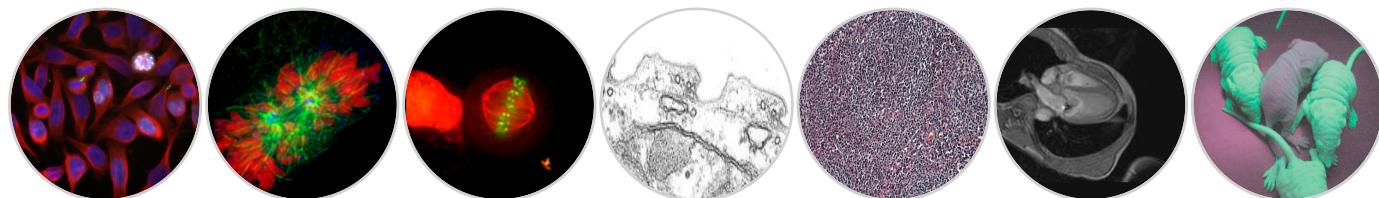
# Reference images



## Biomolecular Resources



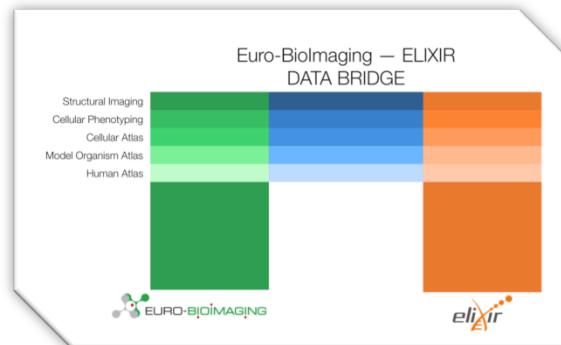
## Imaging Domains



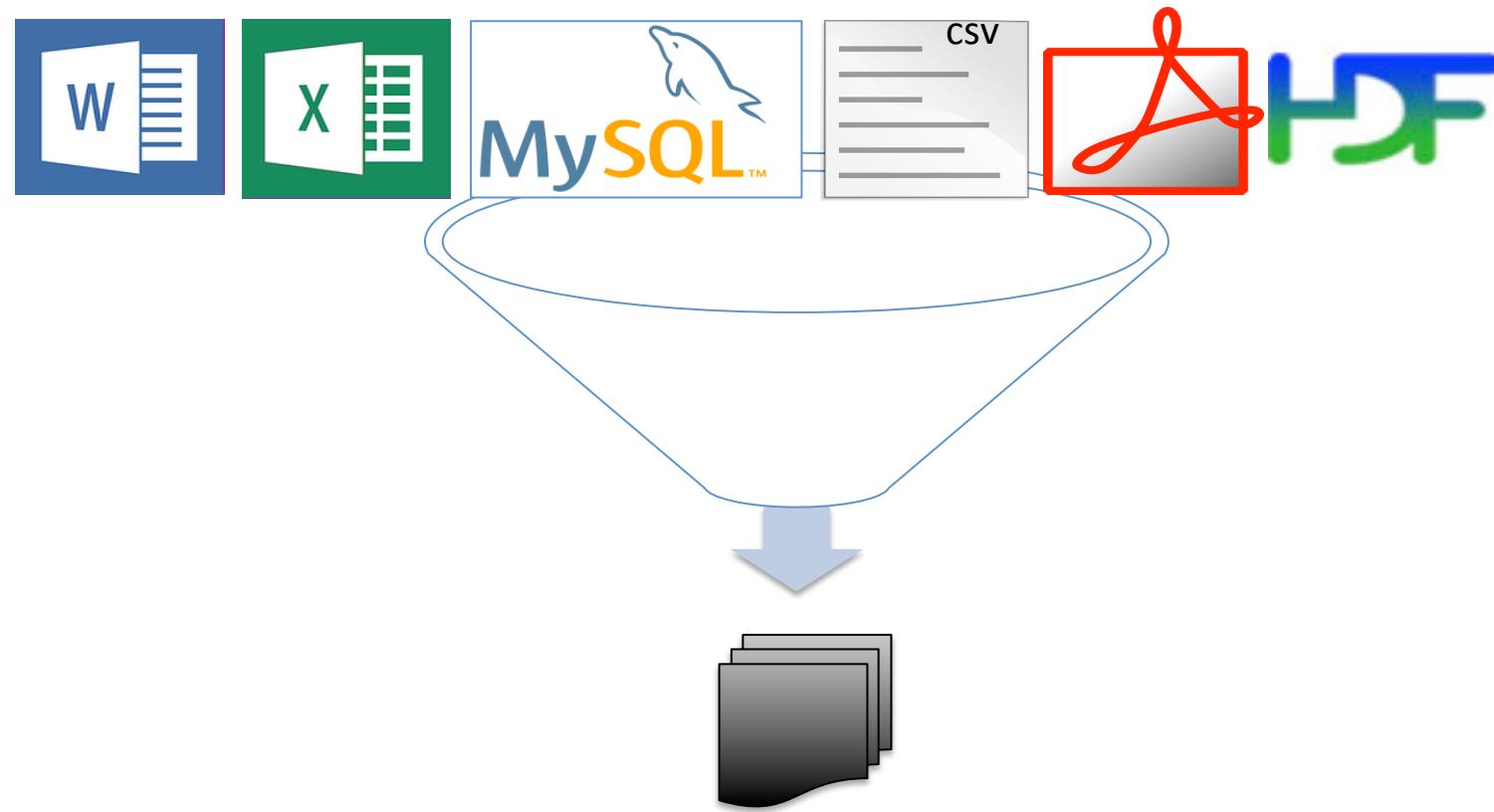
## Controlled Vocabularies



Google: "Euro-BioImaging/Elixir Data Strategy"



# Submitted metadata



*In the style of:* MAGE-TAB **isatab** or, Cellular Phenotype Database

# Vital Stats: IDR @ EBI Embassy



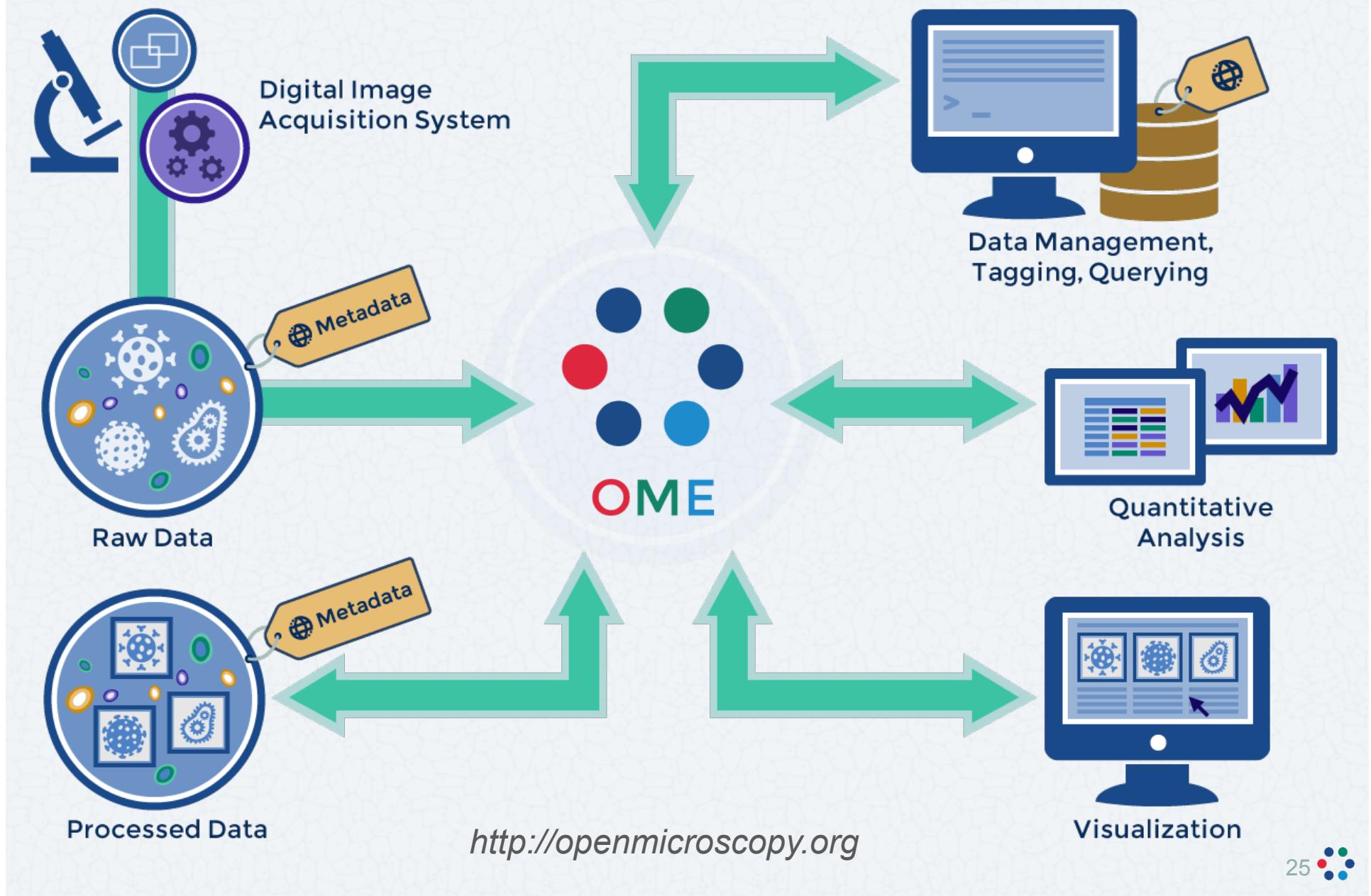
<http://idr.openmicroscopy.org>

Metric	Sept 2017
Image Data Size	42.8 TB
Image files	14.6 M
Datasets	3805
Total Images	2.68 M
Planes	37.3 M
Experiments	1.04M
Genes	19,605
Annotated Images	386 k
Phenotypic Classes	161
External Links	172 k

→~\$25,000/month on AWS←

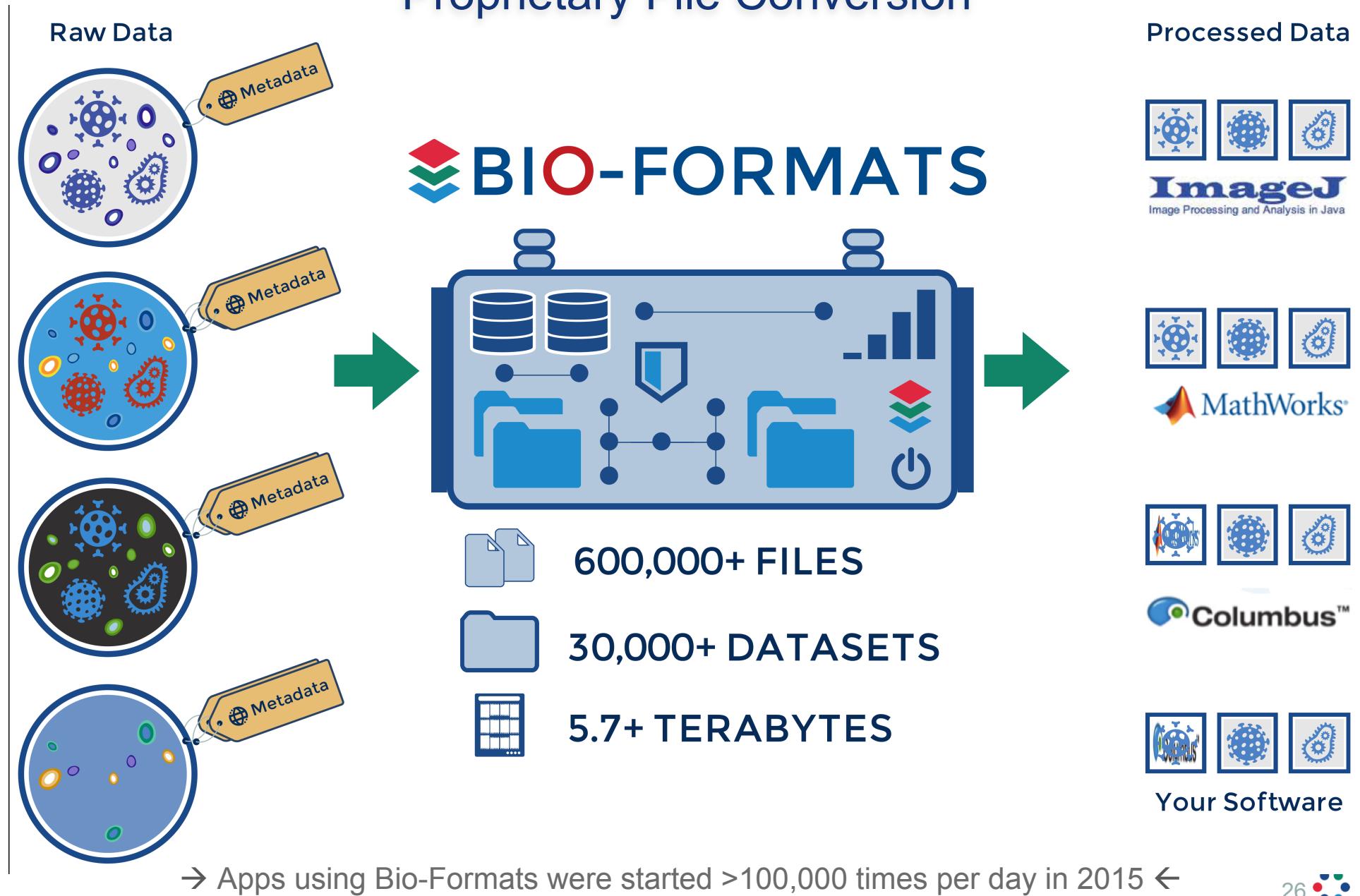
*Williams et al (2017) Nature Methods*

# ...Towards Image Informatics



# BIO-FORMATS:

## Proprietary File Conversion



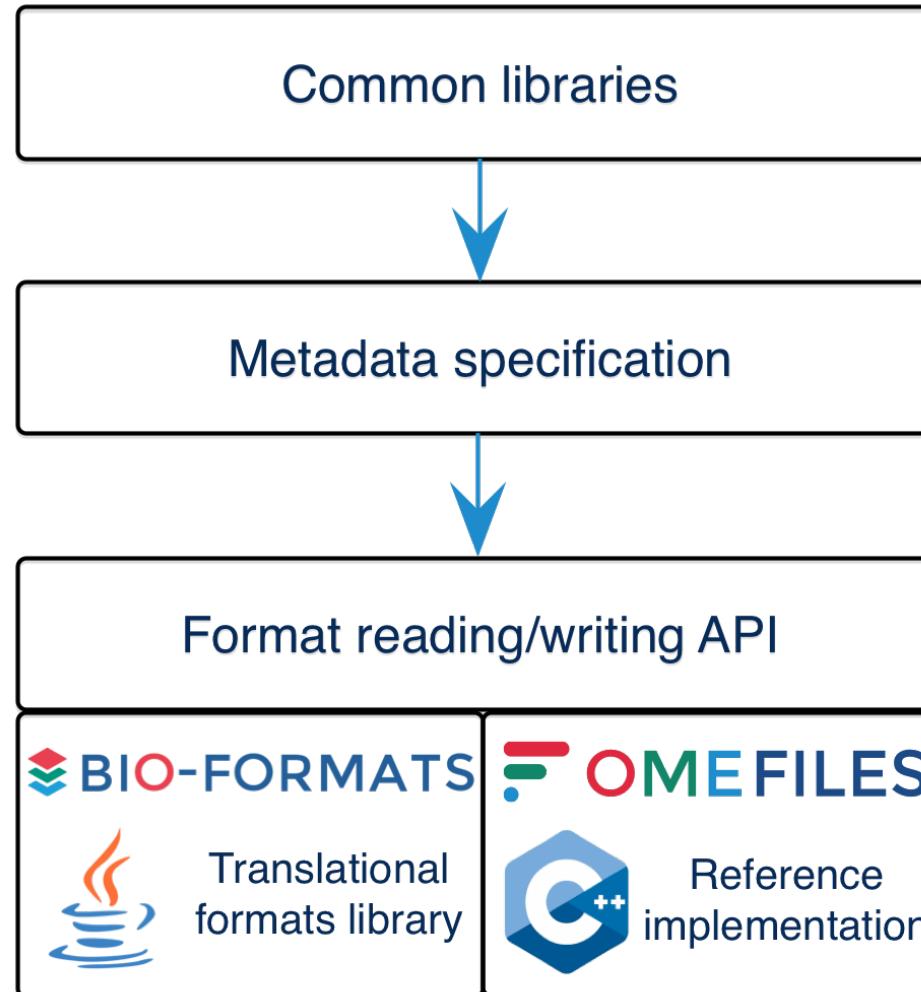
# BIO-FORMATS:

## Opportunities for Academic/Commercial Collaborations

- 3i: Building and maintaining a SlideBook file format reader
- PerkinElmer: Commissioning Glencoe to provide open source Harmony HCS format reader
- ZEISS: Commissioning Glencoe to build open source JPEG-XR decoder

→ See *blog.openmicroscopy.org* for more info

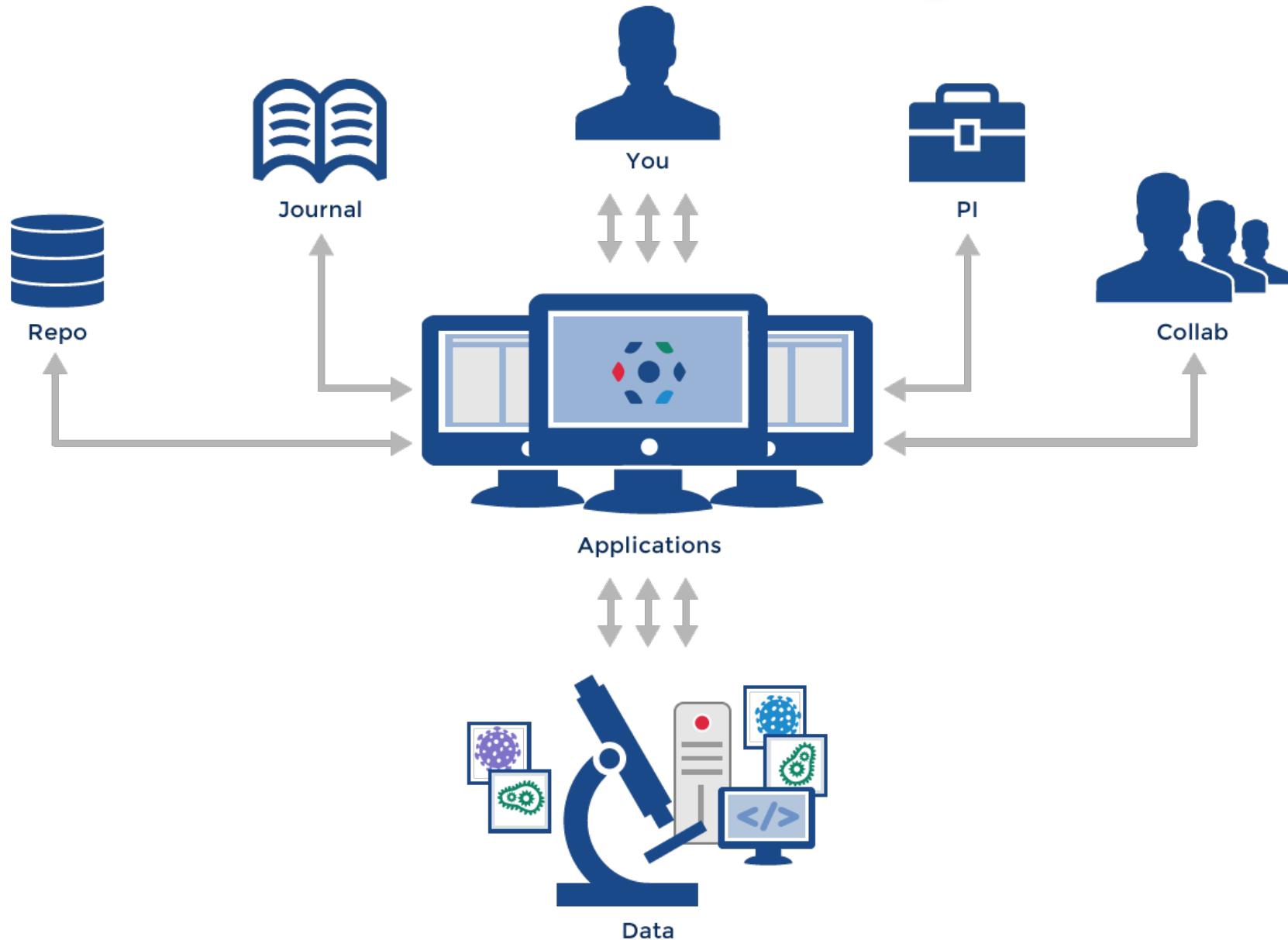
# Software Stack



# The Standard Paradigm

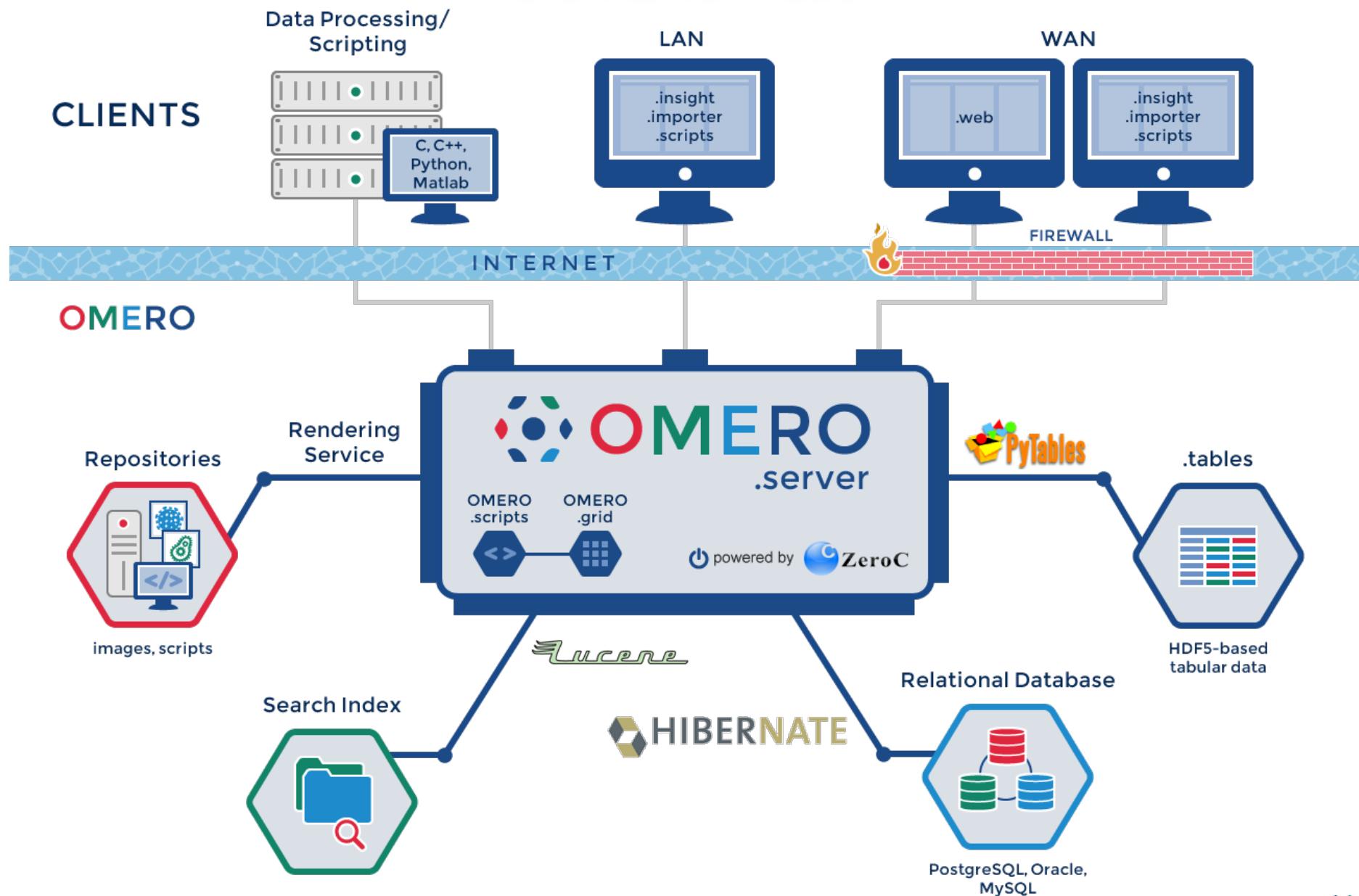


# The “Scientific Data” Paradigm

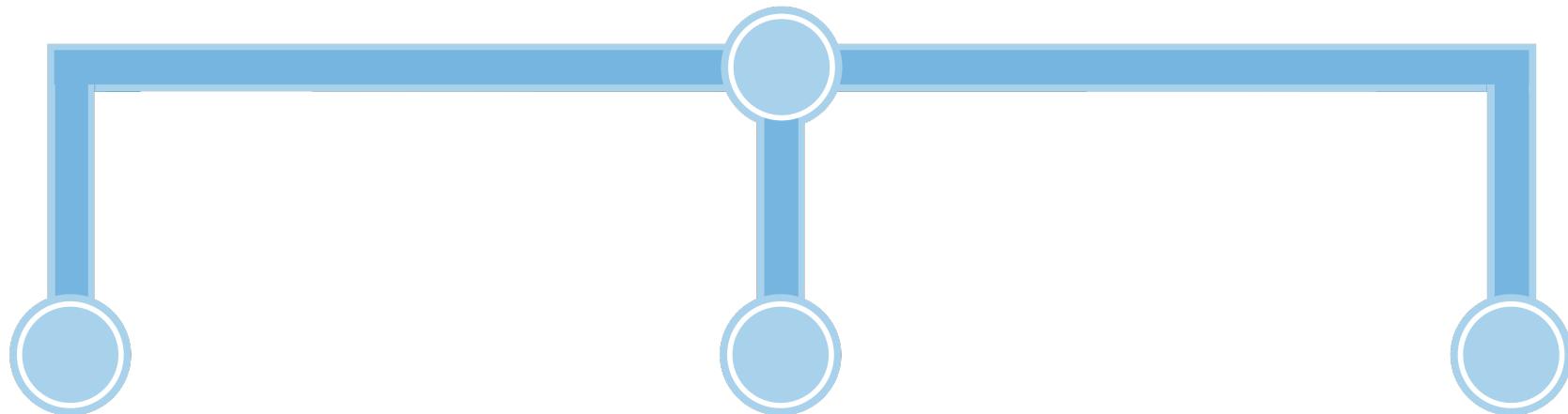


Gray et al, 2005, *Scientific Data Management in the Coming Decade*, Microsoft Research

# The OMERO Platform



# What We Do



**OME-XML**  
**OME-TIFF**  
**OME FILES**

Open  
exchangeable  
file formats  
&  
metadata

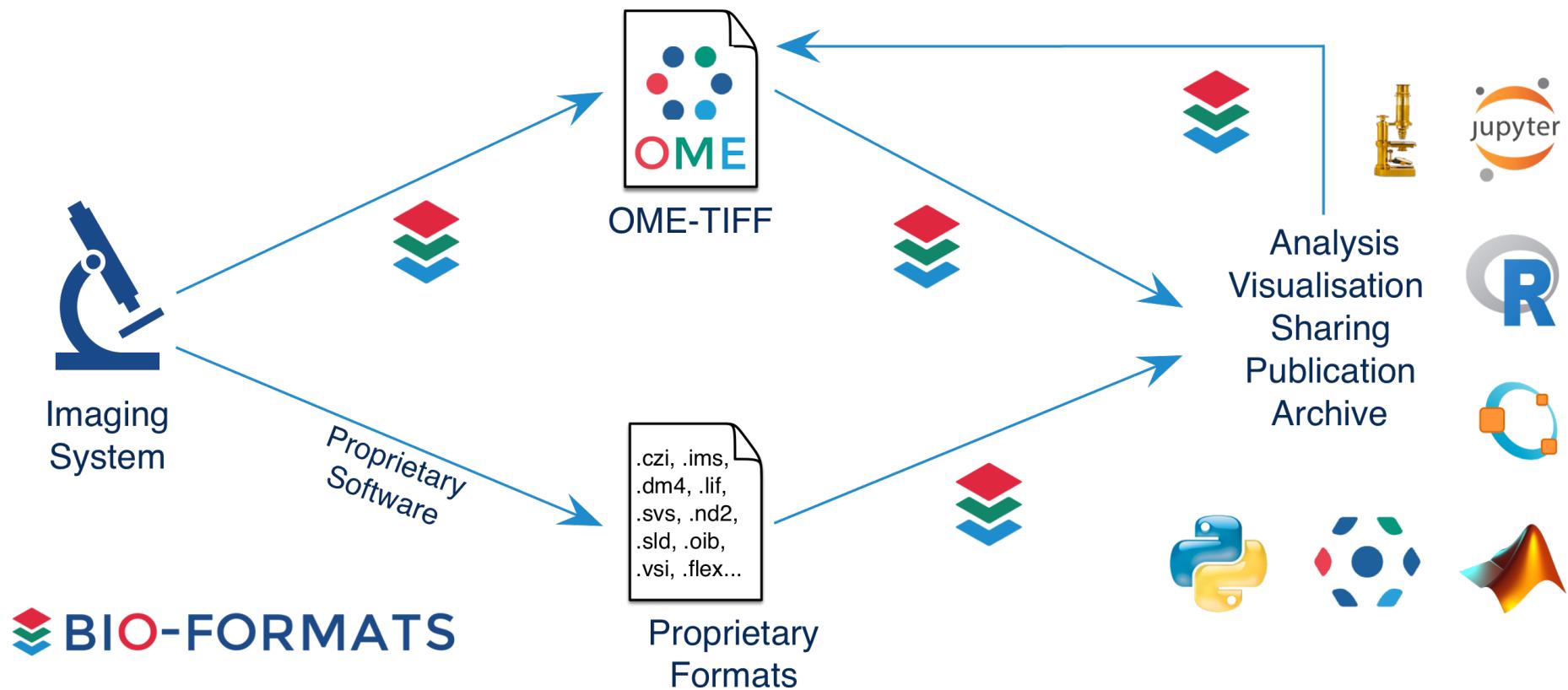
**BIO-FORMATS**

Open  
Proprietary File  
Format  
Translation

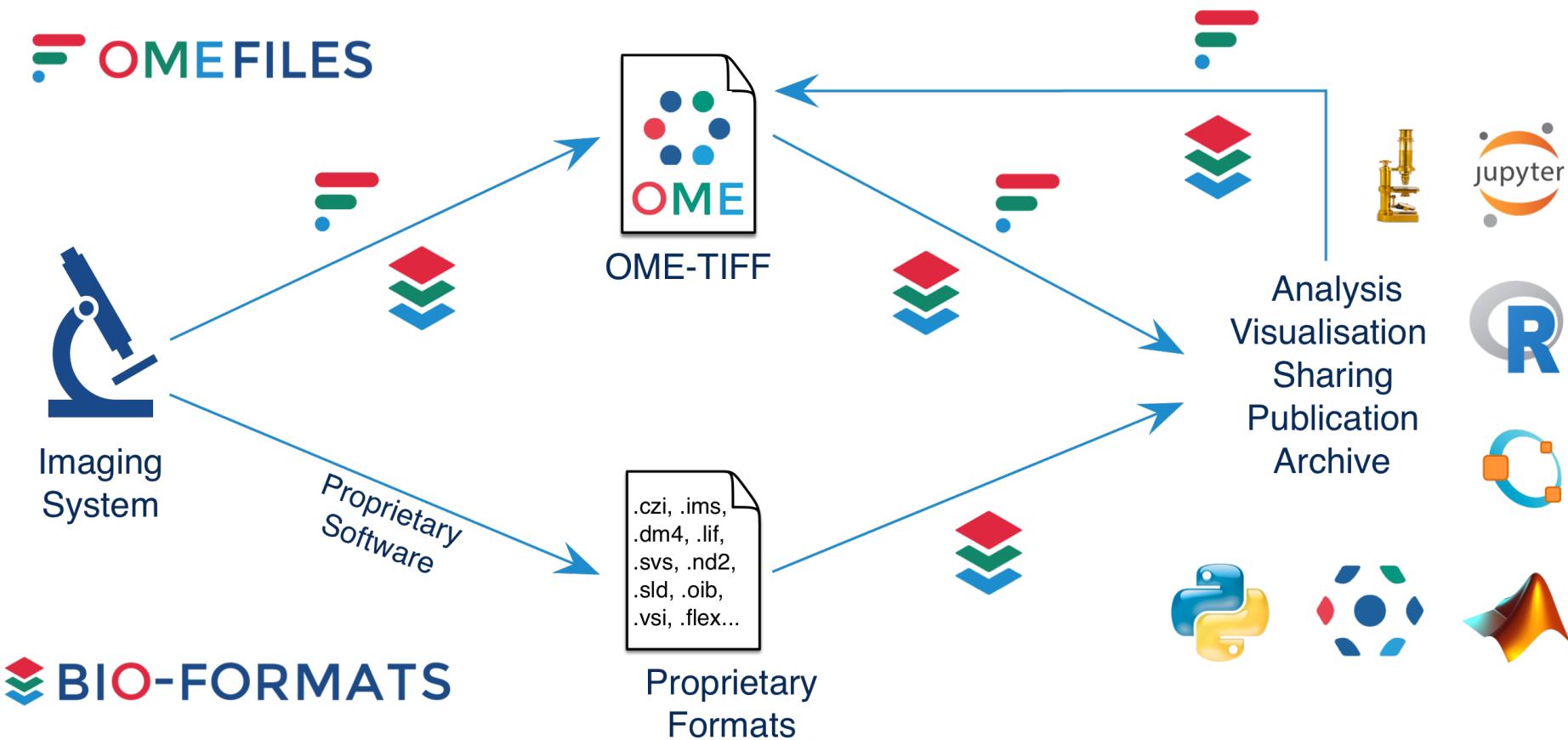
**OMERO**

Open  
Image Data  
Management

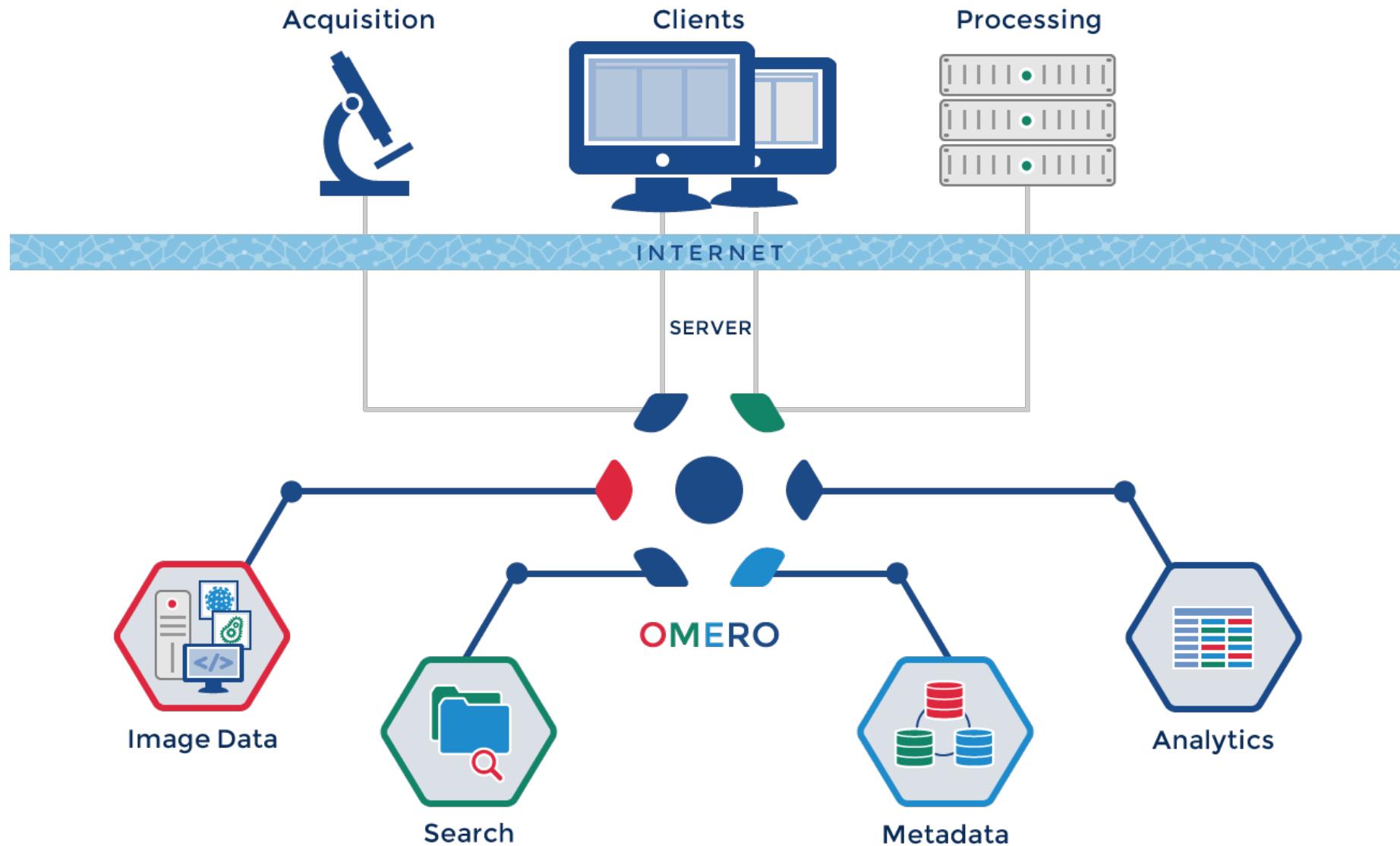
# Bio-Formats



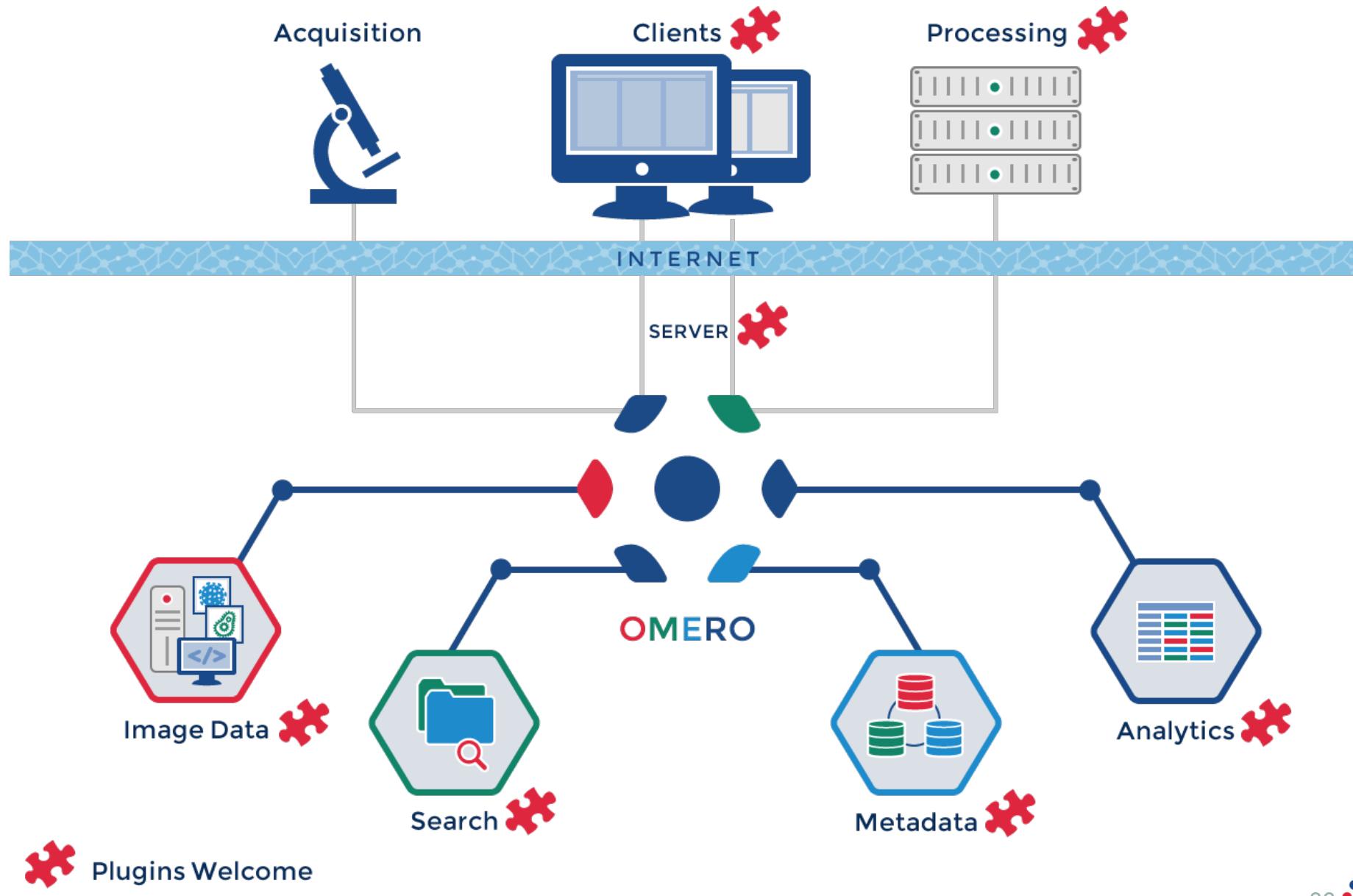
# Bio-Formats



# The OMERO Platform

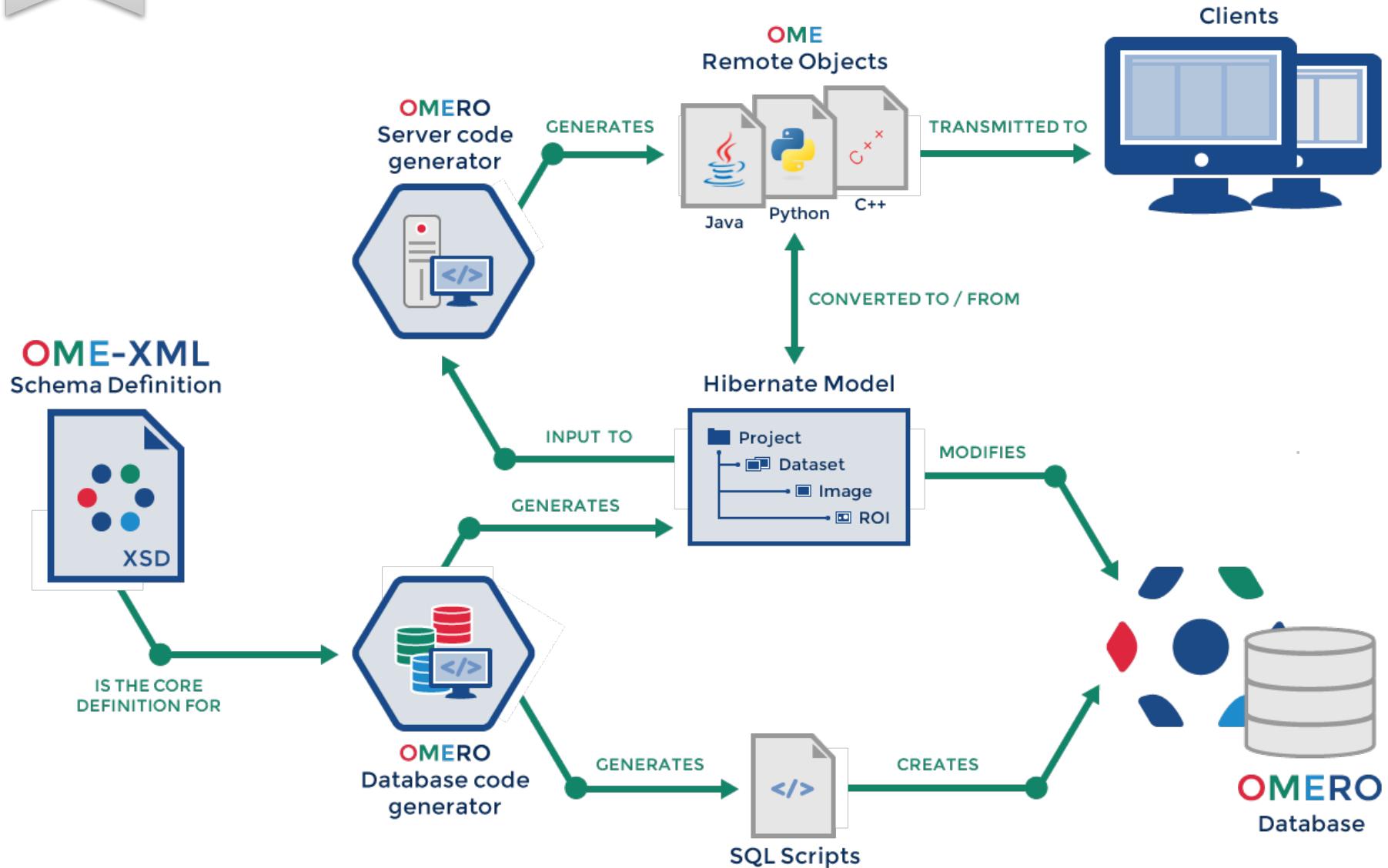


# The Extensible OMERO Platform





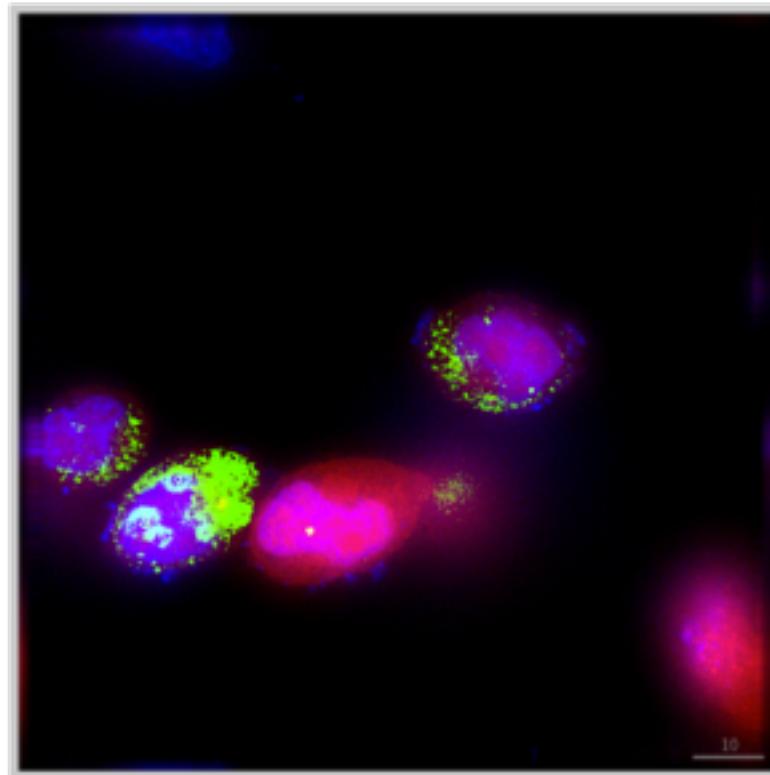
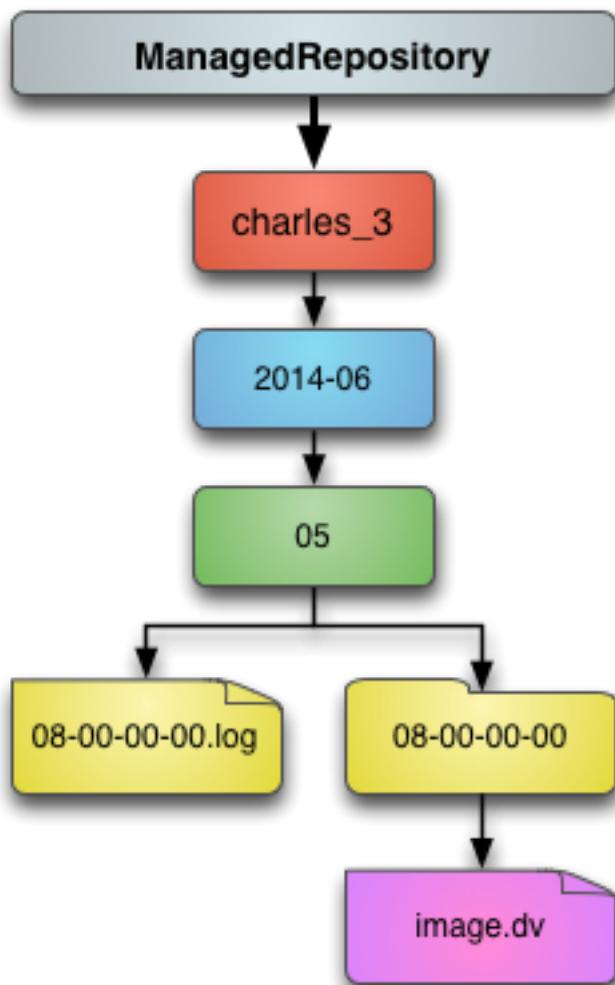
# Leveraging the OME Data Model



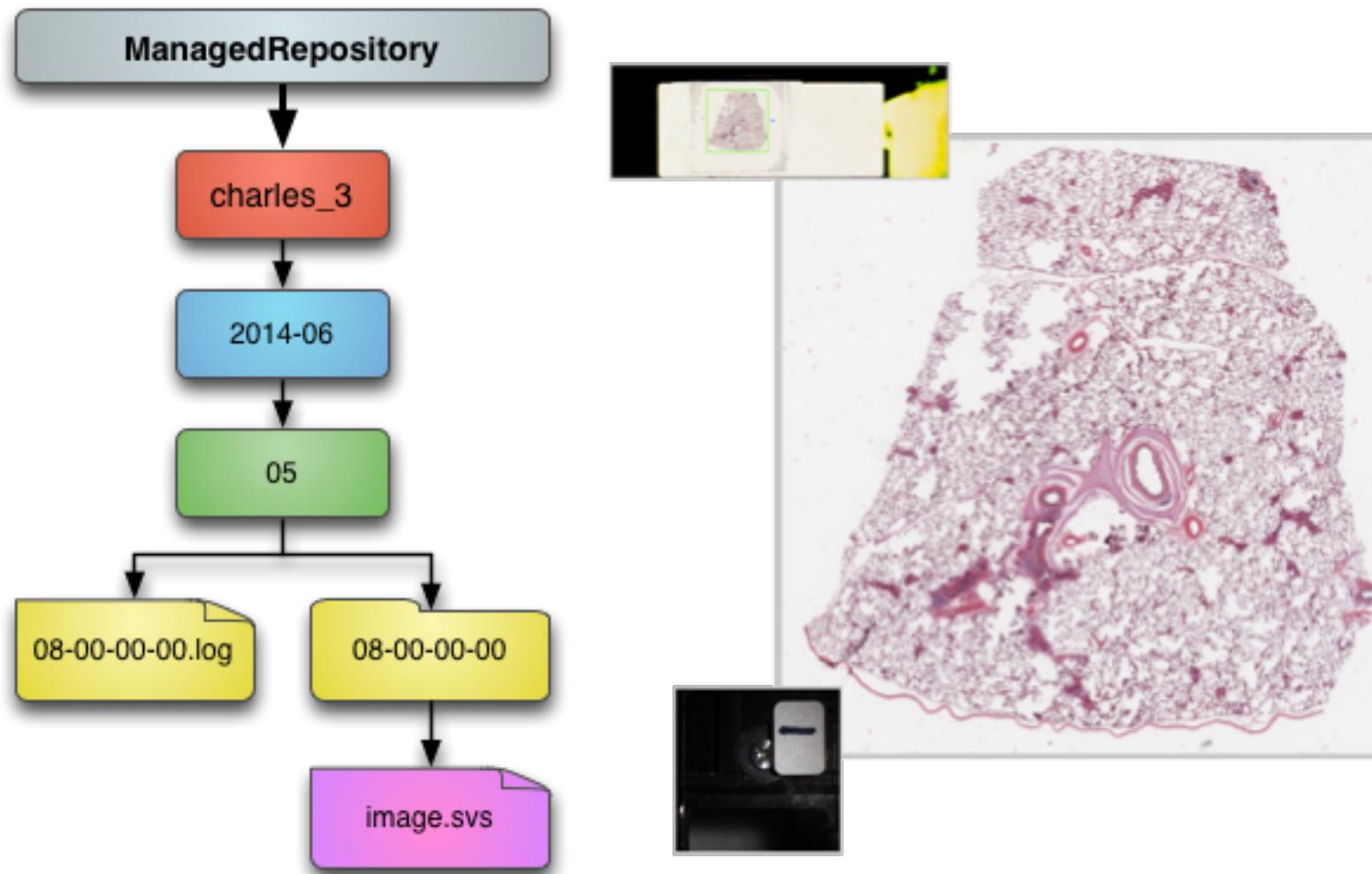
# Security Model

PERMISSIONS	 Read	 Annotate	 Write	 Privacy
	Private	Group-Read	Group-Annotate	Your Group
	Group-Write	Public-Read	Public-Annotate	Anyone
 Private	✓	✓	✓	 You
 Group-Read	✓	✗	✗	
 Group-Annotate	✓	✓	✗	 Your Group
 Group-Write	✓	✓	✓	
 Public-Read	✓	✗	✗	
 Public-Annotate	✓	✓	✗	 Anyone
 Public-Write	✓	✓	✓	

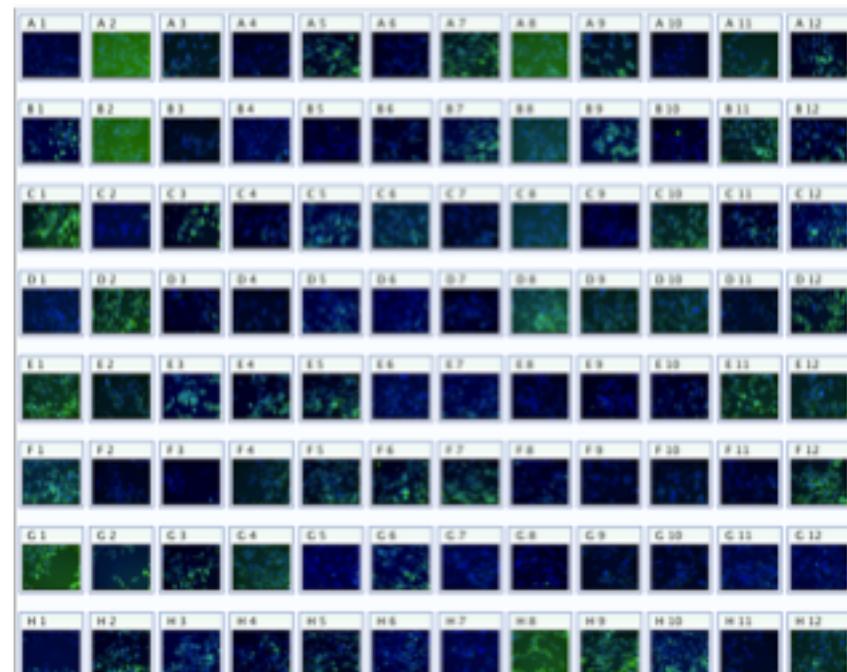
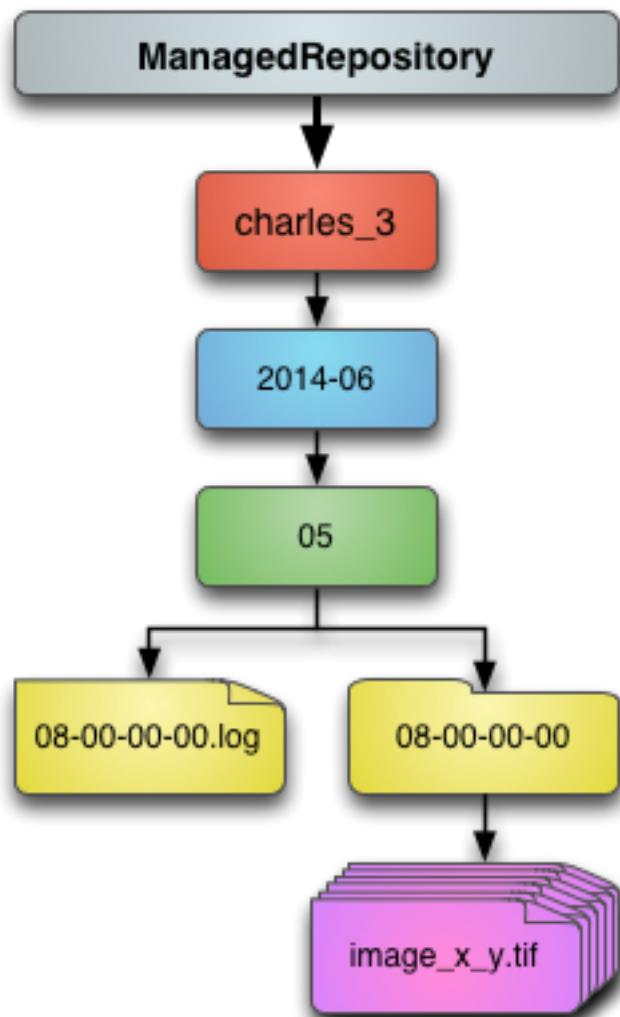
## Filesets: 1-1



## Filesets: 1-N



## Filesets: N-N



# Modulo Dimensions

```
<SA:StructuredAnnotations>
  <SA:XMLAnnotation ID="Annotation:3" Namespace="openmicroscopy.org/omero/dimension/modulo">
    <SA:Value>
      <Modulo namespace="http://www.openmicroscopy.org/Schemas/Additions/2011-09">
        <ModuloAlongZ Type="angle" Unit="degree">
          <Label>45</Label>
          <Label>90</Label>
        </ModuloAlongZ>
        <ModuloAlongT Type="lifetime" TypeDescription="TCSPC" Start="0" Step="2" End="128"/>
        <ModuloAlongC Type="phase" Start="0" Step="1" End="255"/>
      </Modulo>
    </SA:Value>
  </SA:XMLAnnotation>
</SA:StructuredAnnotations>
```