

De-identified Digital Pathology Training Session

Gabriele Campanella

Brandon Veremis

Andrew Deonarine

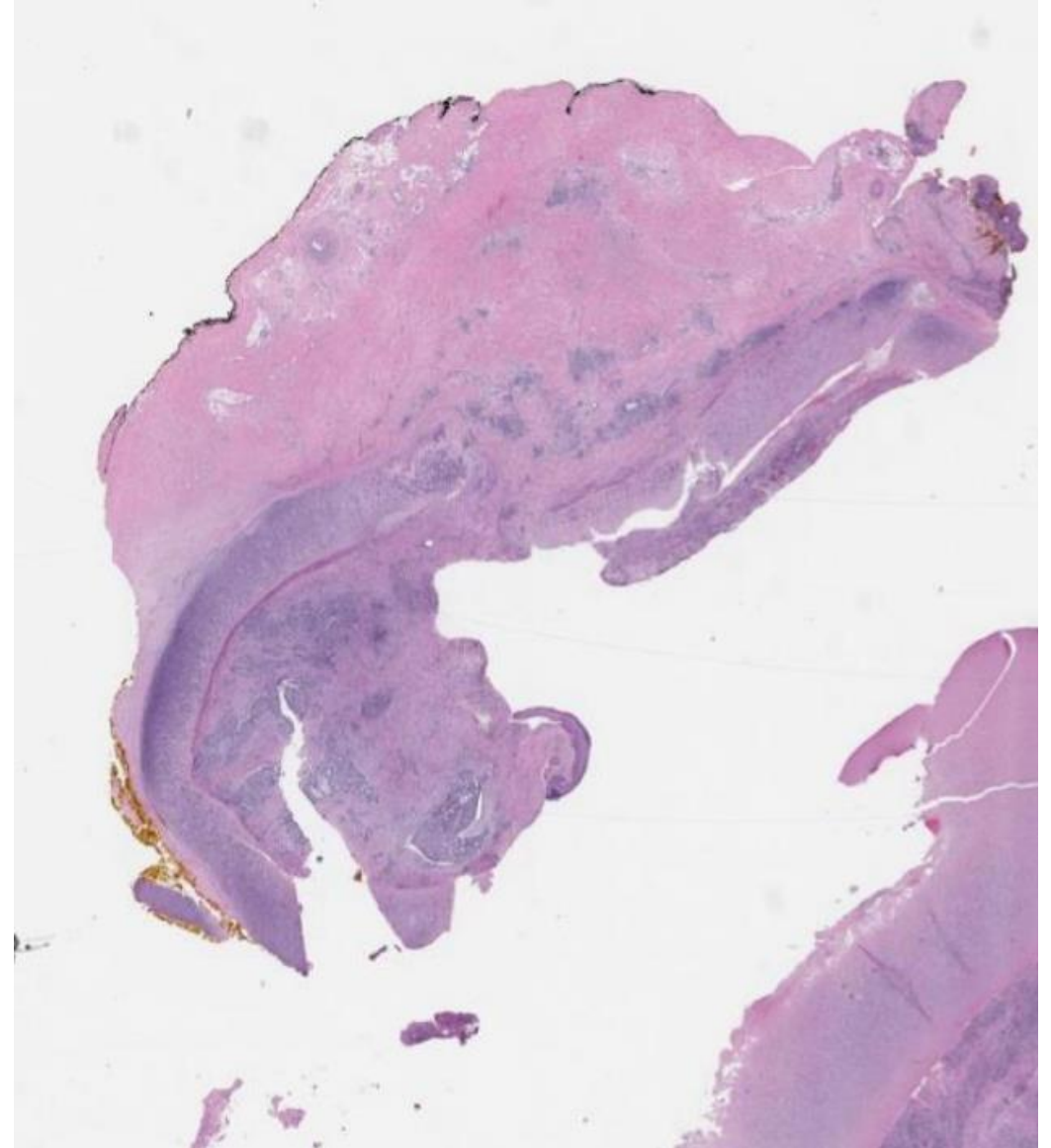
December 2, 2025



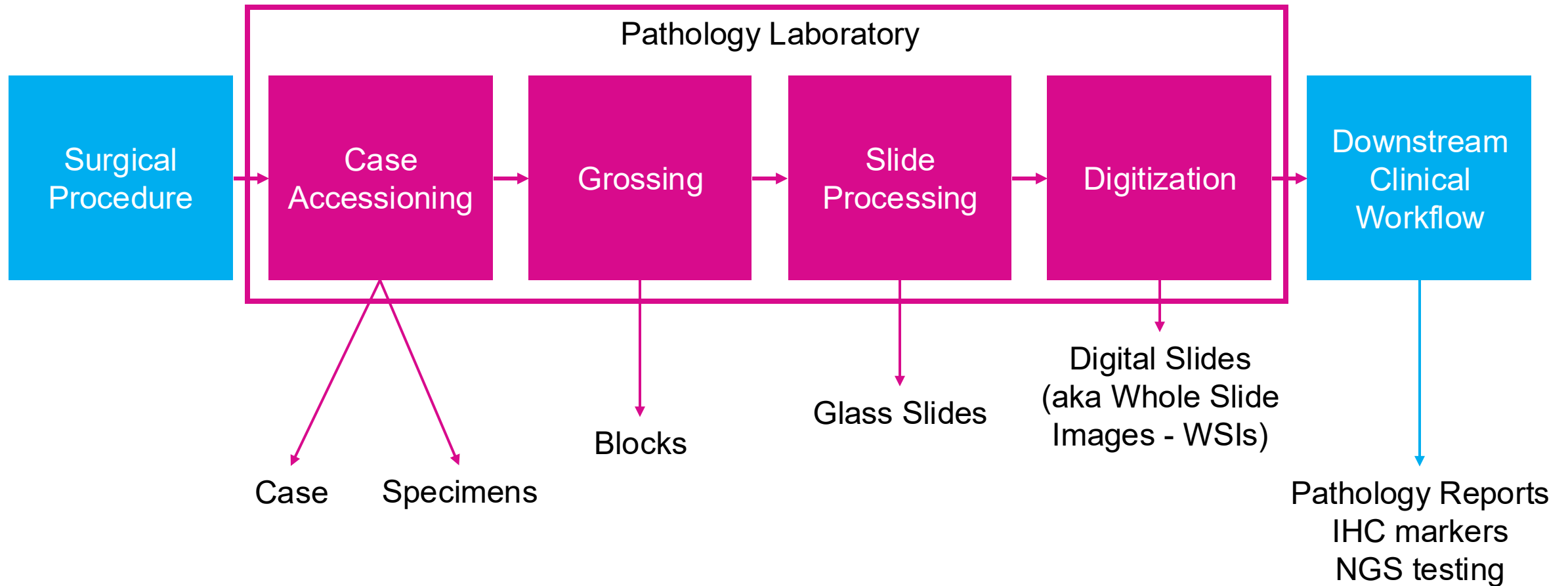
**Icahn School
of Medicine at
Mount
Sinai**

Agenda

- 1) Basic Pathology Asset Structure
- 2) Overview of Pathology Image Data
- 3) How To Interact With Images
- 4) Computational Pathology
- 5) Case Study
- 6) Identify Case/Slide/Cohorts
- 7) Summary



Pathology Lab Workflow and Data



Overview of Pathology Data

The pathology department produces large amounts of physical and digital data.

Data Modalities

- Biological Assets → Stored for a minimum of 20 years per state regulations
 - Tissue blocks
 - Glass Slides
- Digital Assets
 - Pathology Laboratory Information System (LIS) Metadata → Links cases, specimens, blocks, slides to reports. Tracks most aspects of pathology lab
 - Whole Slide Images → Pyramidal .tiff files scanned with Philips WSI scanners at 40x resolution. Available as part of Data Ark
 - Pathology Reports → Free text diagnosis written by the pathologist for the entire case. Can include test results (IHC, FISH, etc.) as free text.

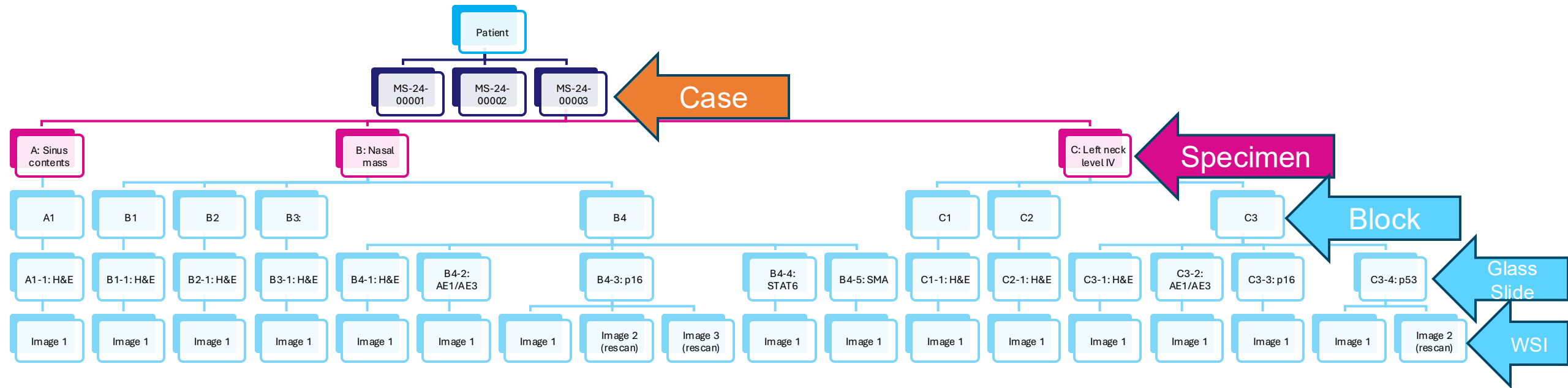
Basic pathology asset structure

One to many to many to many to many to many

- A patient can have many cases
- A case can have many specimens
- A specimen can have many blocks
- A block can have many slides (which can be different stains/orders)
- A slide can have many images (i.e., rescans)
- No designated "best image" for a patient



Pathology Data Hierarchy



Digital Pathology Statistics

Pathology Service (after 2010)

- Slides : 16,452,652
- Cases: 2,462,596
- Patients: 1,402,535

Digital Pathology (prospective since 2020, select retrospective)

- Slides: 3,892,209
- Cases: 618,004
- Patients: 440,169

Data Ark

- Slides: 2,593,332

Overview of Pathology Image Data



Large (easily > 2 GB each)



40x resolution (~0.25 microns/pixel)



Contains all tissue detected by the scanner (usually reliable)



Most images rescanned if QC issues present (e.g., out of focus)



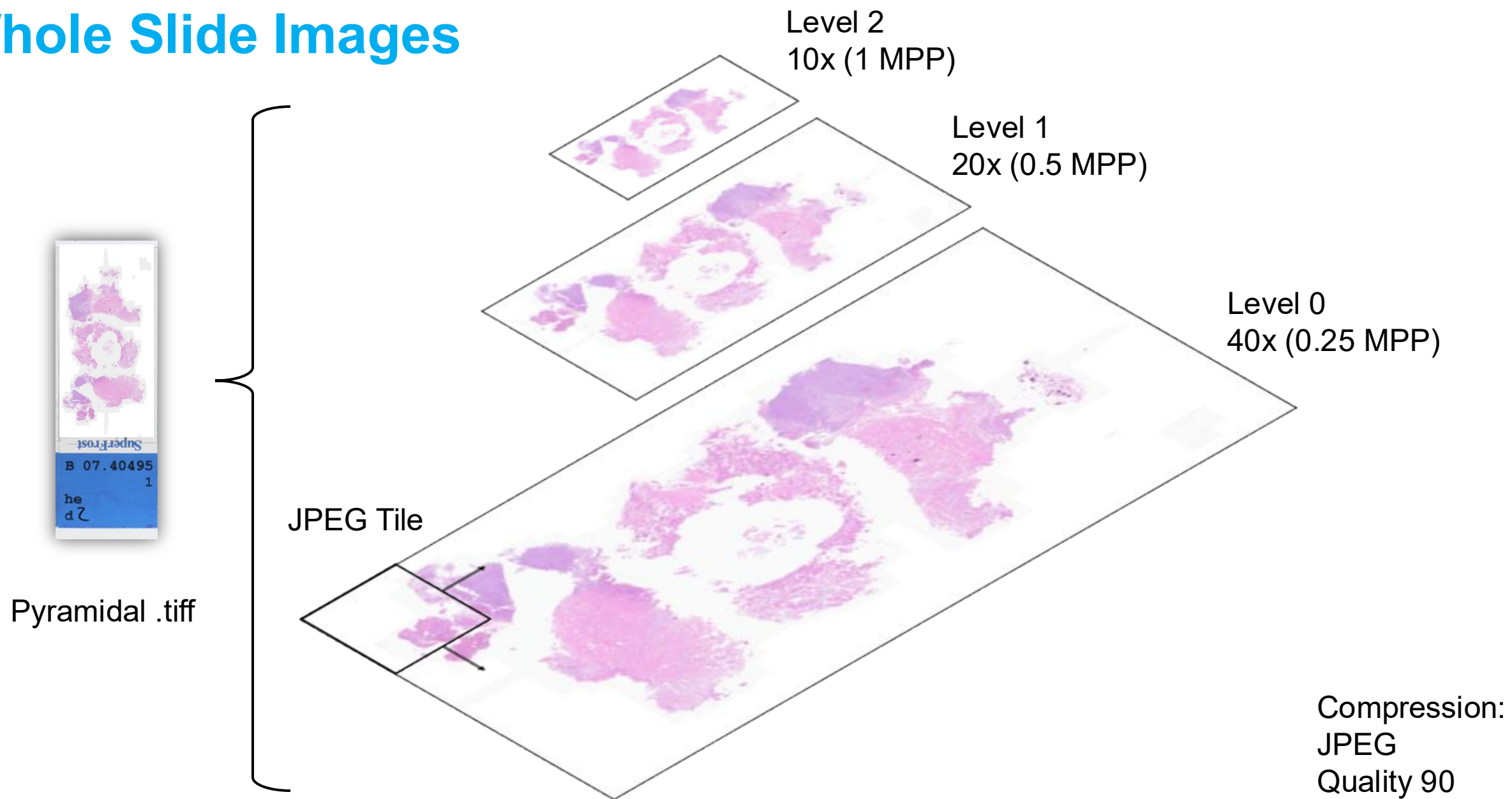
TIFF format



Open-source tools available



Whole Slide Images



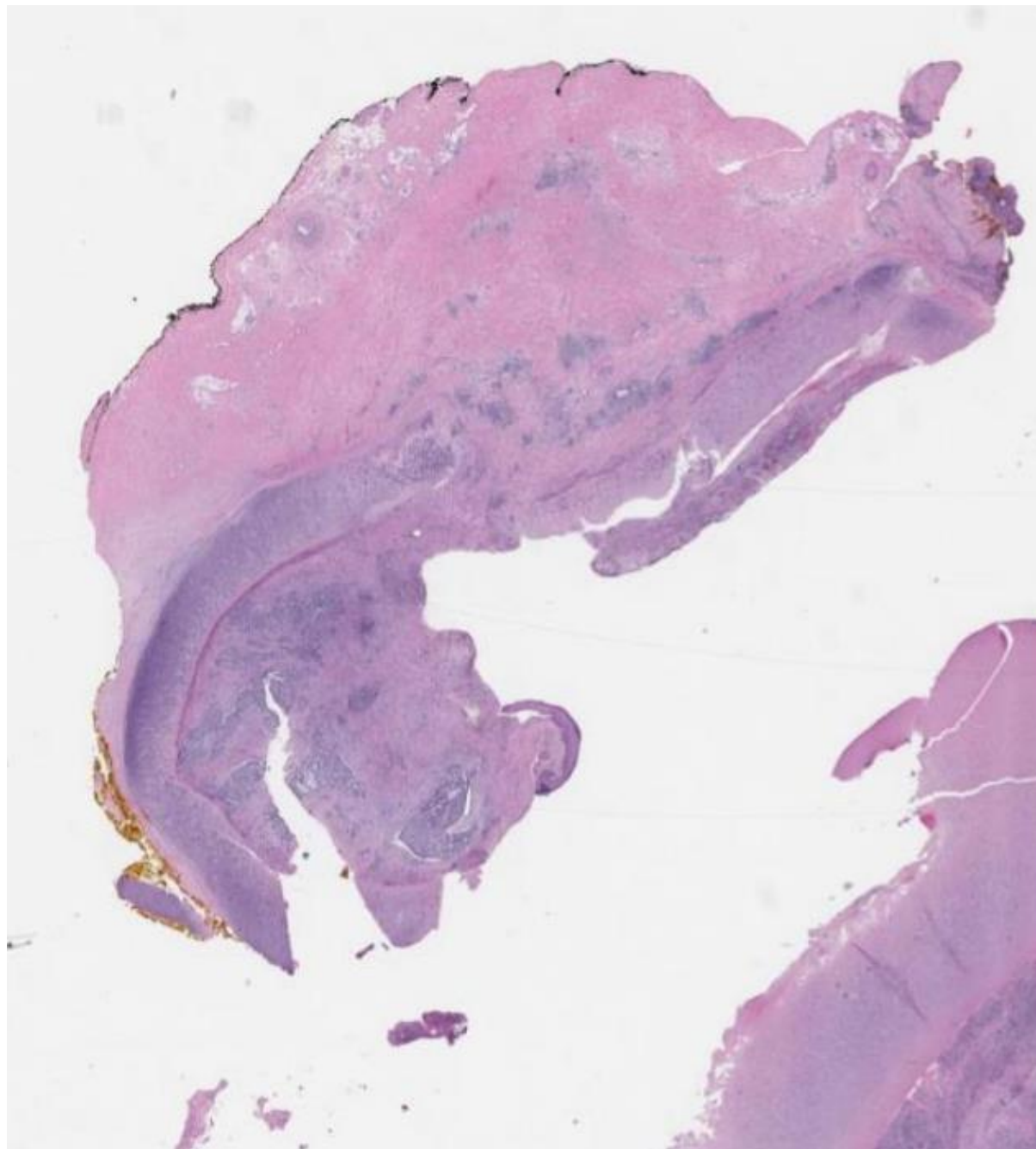
How To Interact With Images

- Files can't be loaded into memory all at once
- Image has to be tiled or downsampled for analysis
- Need specialized packages to interact with the files (python)
 - `openslide`
 - `cucim`
- Needs specialized viewers to view the files
 - QuPath (local)
 - Digital Slide Archive (browser app with backend)

How To Interact With Images

```
import openslide
slide = openslide.OpenSlide('path/to/slide.tiff')
base_mpp = float(slide.properties[openslide.PROPERTY_NAME_MPP_X])
thumb = slide.get_thumbnail((1000,1000))
downsamples = slide.level_downsamples
tile40x = slide.read_region((14356, 24674), 0, (256,256)).convert('RGB')
tile20x = slide.read_region((14356, 24674), 1, (256,256)).convert('RGB')
```

<https://openslide.org/api/python/>



Quality Control

- Detect blurred regions
- Detect tissue folds
- Stain normalization

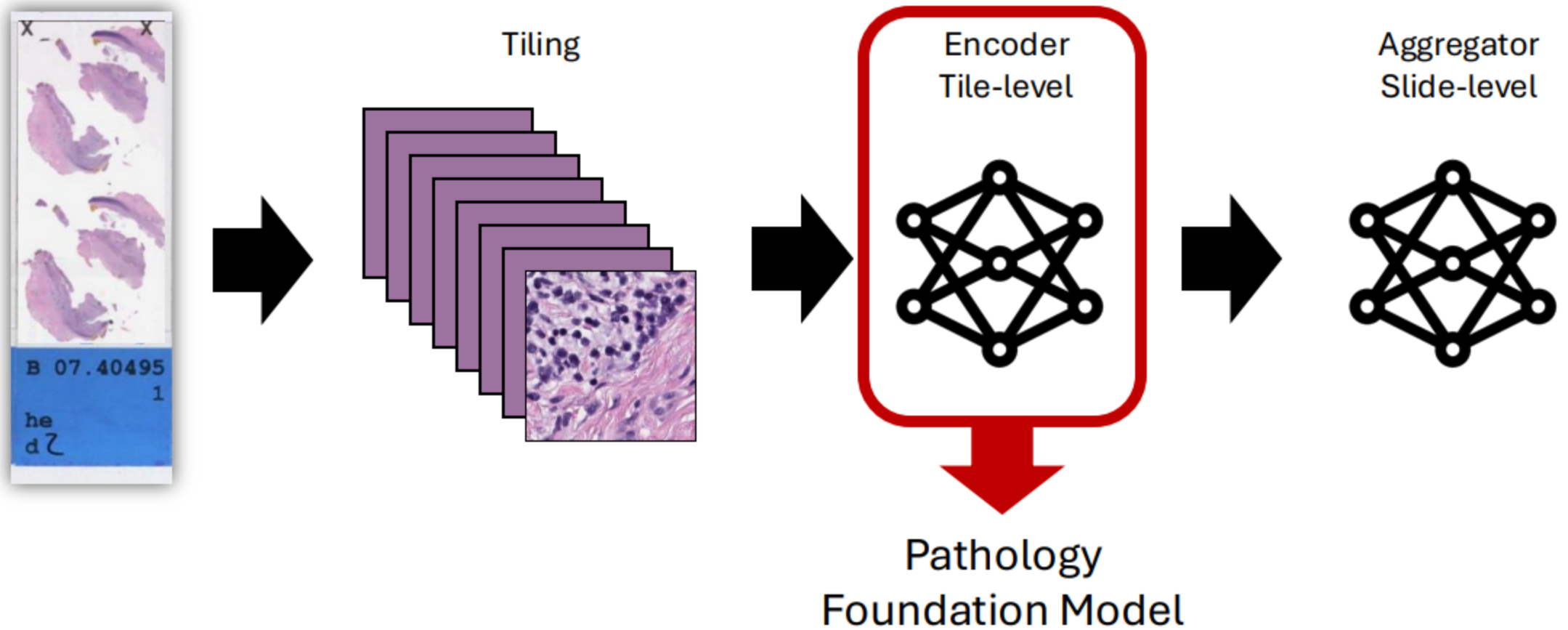
Diagnostic Reporting

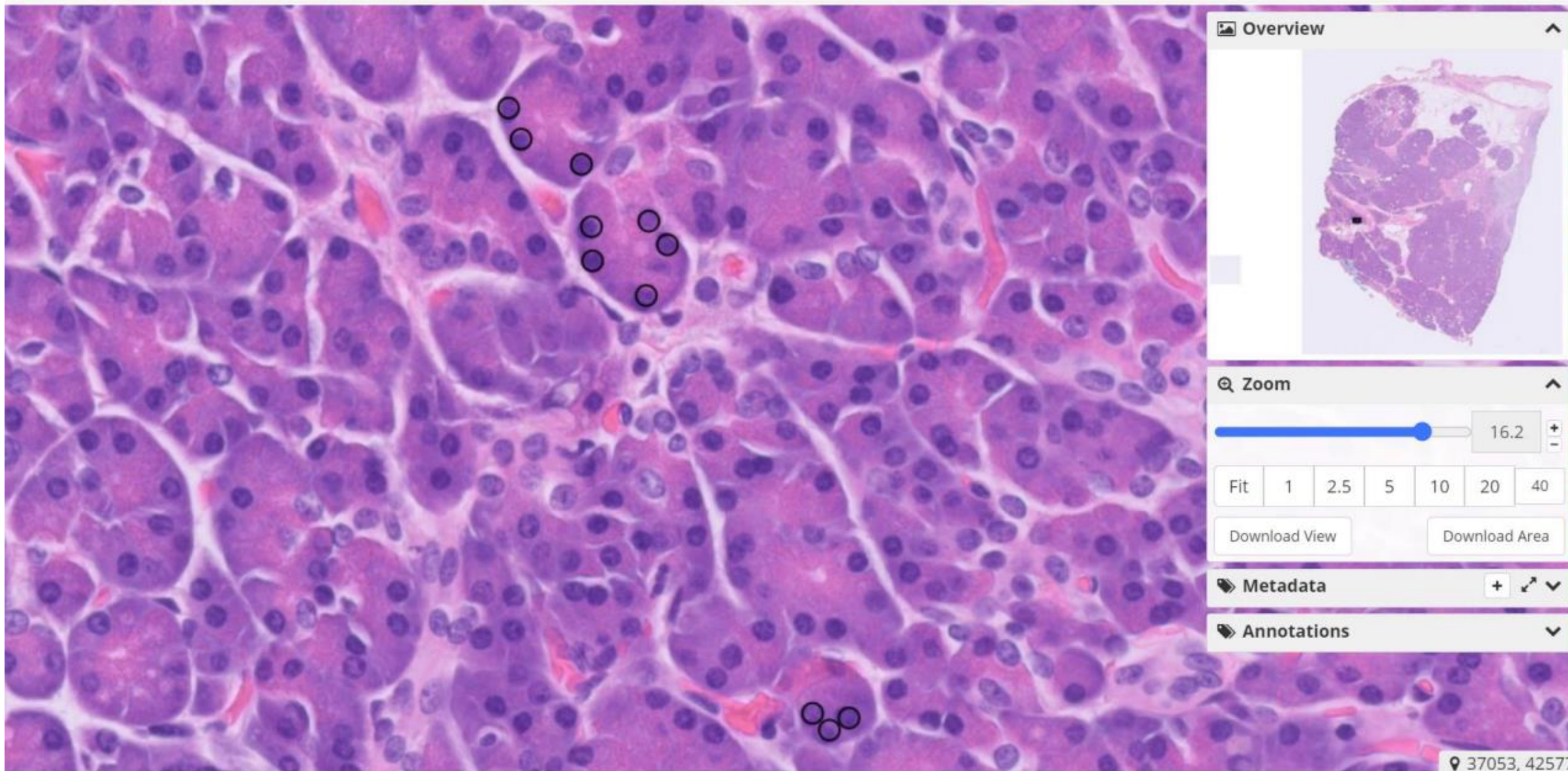
- Tumor detection
- Tumor segmentation
- Tumor staging
- Cell detection
- Slide retrieval

Prognostication

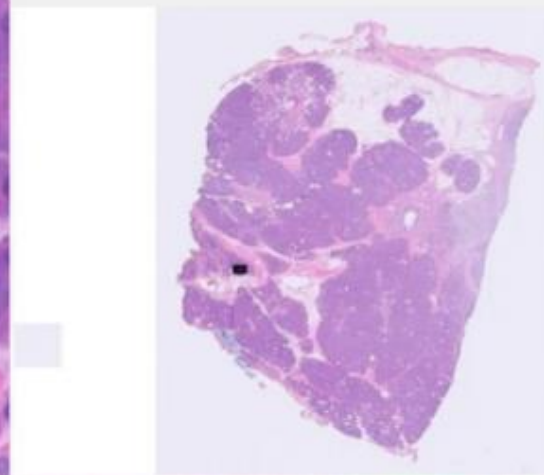
- Biomarker prediction
- Response prediction
- Survival analysis

Computational Pathology

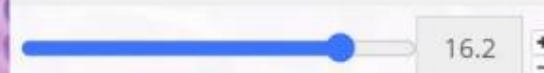




Overview



Zoom



Fit

1

2.5

5

10

20

40

Download View

Download Area

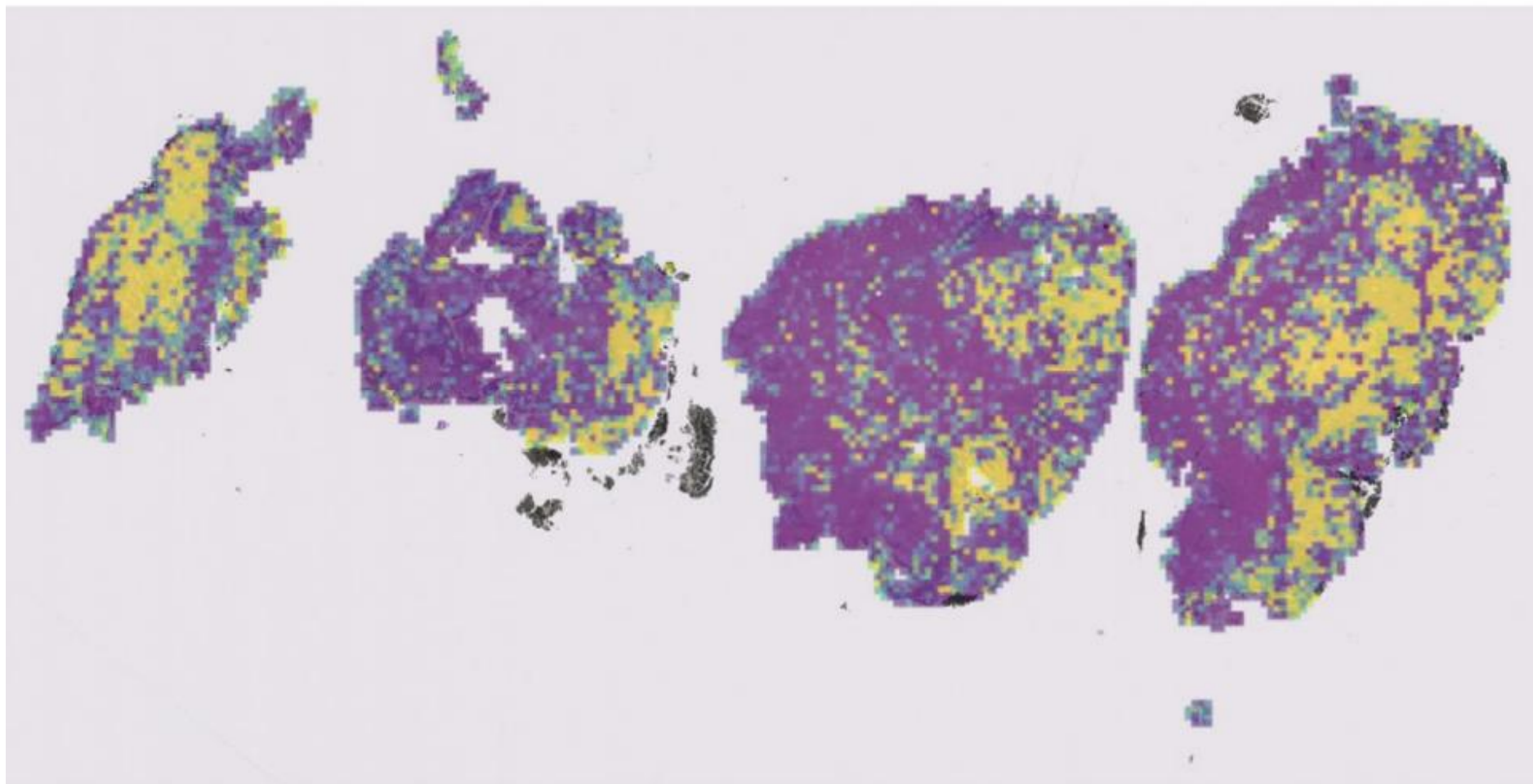
Metadata



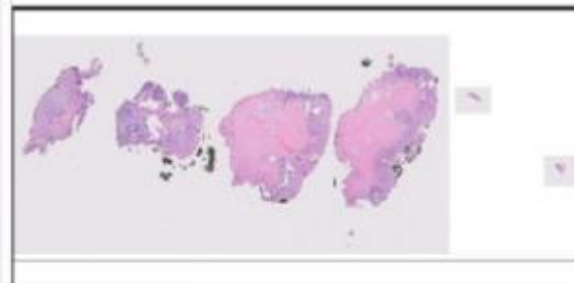
Annotations



37053, 42577



Overview ^



Zoom ^



Metadata + ↗ ▾

Annotations ▾

Case Study

Goal: Obtain a diagnostic biopsy slide cohort of breast cancer patients with invasive ductal carcinoma and their receptor status.

Tasks:

- Identify the right cases
- Identify representative slide
- *Obtain diagnosis, filter by IDC*
- *Obtain hormone receptor status (ER/PR/HER2)*

Identify Case

- Manual review
- Rely on pathology Sign-out knowledge
 - Search breast cases where report contains “biopsy” & “ductal carcinoma”
- Cons
 - Will capture DCIS
 - May lose cases where biopsy is not explicitly mentioned
 - Typos
 - Free text fallacy

Identify Case

- Rely on non-pathology data sources
 - EPIC, cancer registry, etc.
 - Match with collection dates of specimens
- Cons
 - Will not capture biopsies done in outpatient clinic or consults
 - Potential inconsistent MRNs

Identify Slide

- Preface: pathology report gives specimen level information, not slide level
 - Manual slide review
 - Rely on specialty specific knowledge
 - For breast cancer, only tumor slides receive ER/PR/HER2 testing
- Rely on synoptic reporting
 - Newer synoptic reports include “best” tumor block
 - Cons
 - Implemented only recently
 - No synoptic for biopsies
- AI – vision model

Pathologist collaboration is essential!

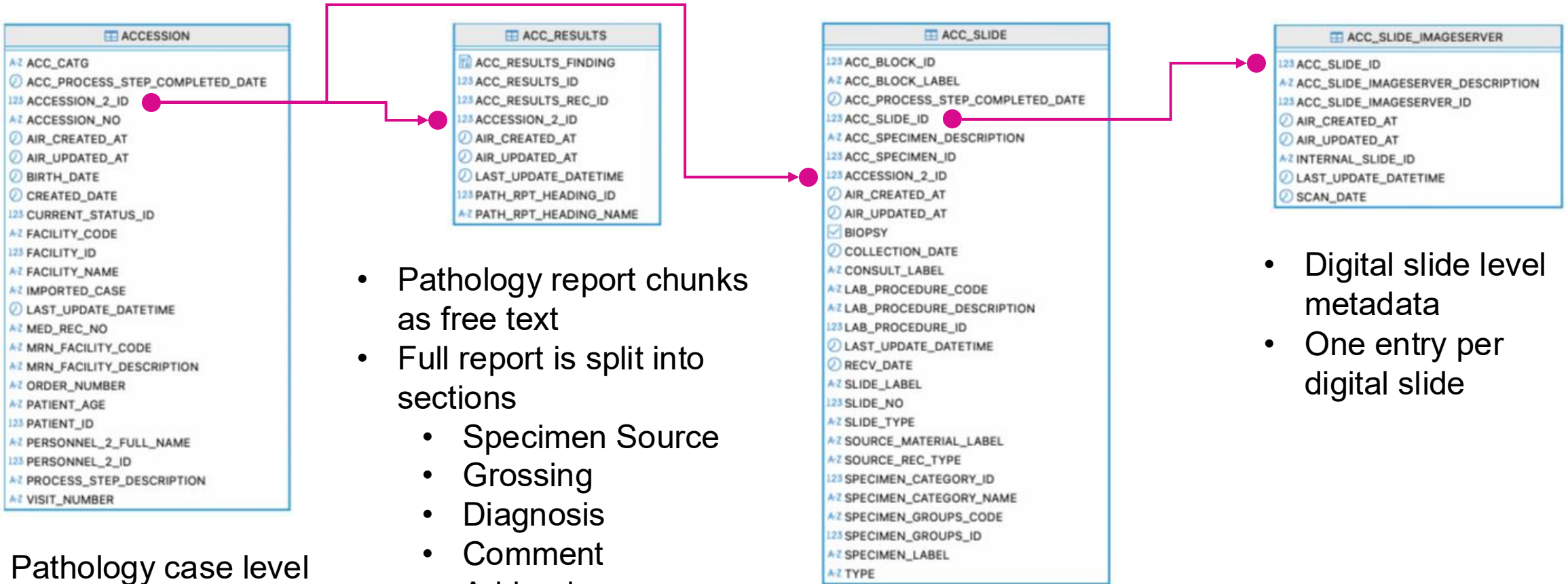
Using Leaf to Locate Digital Pathology Images on Data Ark

- Example of file linkages in Leaf (to the file in Data Ark)
- These are linked at the patient level, so you will get a list of files (working on refining this linkage)
- We are in the process of incorporating this into AIRMS

The screenshot displays the Leaflet application interface. At the top, there's a navigation bar with the 'leaf' logo, a search bar, and a 'New Query' button. Below this, a header indicates '14,951 patients'. The main content area shows a table of patient data with columns: Person ID, Patient of, Address, Address State, Age, Ethnicity, Gender, Language, Marital Status, Race, Religion, and Pathology Slides File Paths. The table lists 8 patients, each with a 'View details' link. The bottom of the screenshot shows a blue bar with file paths, which are partially obscured by a redacted area.

| Person ID | Patient of | Address Postal Code | Address State | Age | Ethnicity | Gender | Language | Marital Status | Race | Religion | Pathology Slides File Paths |
|------------|---------------------------|---------------------|---------------|-----|------------------------|--------|----------|----------------|-------------|----------|---|
| 504990694 | Mount Sinai Health System | 100 | NY | 79 | Not Hispanic or Latino | MALE | Unknown | Unknown | Chinese | Unknown | /sc/arion/projects/data-ark/digital_pathology_slides/e9c/e9c21f63-7bca-42bf-a9a3-22169e8e6c0e.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/051/05191f89-fcf8-43a5-8126-4b65a8b328.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/013/01300fde-0f8d-4f95-81c9-4651f8f8a0ba.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/220/2202fd96f-a586-4a33-9ae1-4633-9ae1-9f16387b0d69.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/c4d/c4d4a0a5-eef2-4d0b-9354-c21b704cf5.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/4b5/4b5dc7f0-7787-4260-919e-9b2a04acefa.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/54d/54d61389-3d3a-412b-92c7-2eccc48bc909.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/dcc/dcc158e7-2ff6-4fa6-b51b-8d641af9-180b-42de-aaf0.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/a84/a84f2892-b177-4b65-b63a-4f1bdc6bd5a8.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/e74/e748067a-f7f5-4741-9eec-8748b45756ca.tiff |
| 986919534 | Mount Sinai Health System | 100 | NY | 76 | Not Hispanic or Latino | FEMALE | Unknown | Unknown | Japanese | Unknown | /sc/arion/projects/data-ark/digital_pathology_slides/e9c/e9c21f63-7bca-42bf-a9a3-22169e8e6c0e.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/051/05191f89-fcf8-43a5-8126-4b65a8b328.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/013/01300fde-0f8d-4f95-81c9-4651f8f8a0ba.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/220/2202fd96f-a586-4a33-9ae1-4633-9ae1-9f16387b0d69.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/c4d/c4d4a0a5-eef2-4d0b-9354-c21b704cf5.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/4b5/4b5dc7f0-7787-4260-919e-9b2a04acefa.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/54d/54d61389-3d3a-412b-92c7-2eccc48bc909.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/dcc/dcc158e7-2ff6-4fa6-b51b-8d641af9-180b-42de-aaf0.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/a84/a84f2892-b177-4b65-b63a-4f1bdc6bd5a8.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/e74/e748067a-f7f5-4741-9eec-8748b45756ca.tiff |
| 971929531 | Mount Sinai Health System | 100 | NY | 76 | Not Hispanic or Latino | MALE | Unknown | Unknown | Pakistani | Unknown | /sc/arion/projects/data-ark/digital_pathology_slides/e9c/e9c21f63-7bca-42bf-a9a3-22169e8e6c0e.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/051/05191f89-fcf8-43a5-8126-4b65a8b328.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/013/01300fde-0f8d-4f95-81c9-4651f8f8a0ba.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/220/2202fd96f-a586-4a33-9ae1-4633-9ae1-9f16387b0d69.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/c4d/c4d4a0a5-eef2-4d0b-9354-c21b704cf5.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/4b5/4b5dc7f0-7787-4260-919e-9b2a04acefa.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/54d/54d61389-3d3a-412b-92c7-2eccc48bc909.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/dcc/dcc158e7-2ff6-4fa6-b51b-8d641af9-180b-42de-aaf0.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/a84/a84f2892-b177-4b65-b63a-4f1bdc6bd5a8.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/e74/e748067a-f7f5-4741-9eec-8748b45756ca.tiff |
| 530723726 | Mount Sinai Health System | 111 | NY | 70 | Not Hispanic or Latino | FEMALE | Unknown | Unknown | Bangladeshi | Unknown | /sc/arion/projects/data-ark/digital_pathology_slides/e9c/e9c21f63-7bca-42bf-a9a3-22169e8e6c0e.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/051/05191f89-fcf8-43a5-8126-4b65a8b328.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/013/01300fde-0f8d-4f95-81c9-4651f8f8a0ba.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/220/2202fd96f-a586-4a33-9ae1-4633-9ae1-9f16387b0d69.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/c4d/c4d4a0a5-eef2-4d0b-9354-c21b704cf5.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/4b5/4b5dc7f0-7787-4260-919e-9b2a04acefa.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/54d/54d61389-3d3a-412b-92c7-2eccc48bc909.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/dcc/dcc158e7-2ff6-4fa6-b51b-8d641af9-180b-42de-aaf0.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/a84/a84f2892-b177-4b65-b63a-4f1bdc6bd5a8.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/e74/e748067a-f7f5-4741-9eec-8748b45756ca.tiff |
| 1270071962 | Mount Sinai Health System | 100 | NY | 80 | Not Hispanic or Latino | FEMALE | Unknown | Unknown | Korean | Unknown | /sc/arion/projects/data-ark/digital_pathology_slides/e9c/e9c21f63-7bca-42bf-a9a3-22169e8e6c0e.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/051/05191f89-fcf8-43a5-8126-4b65a8b328.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/013/01300fde-0f8d-4f95-81c9-4651f8f8a0ba.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/220/2202fd96f-a586-4a33-9ae1-4633-9ae1-9f16387b0d69.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/c4d/c4d4a0a5-eef2-4d0b-9354-c21b704cf5.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/4b5/4b5dc7f0-7787-4260-919e-9b2a04acefa.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/54d/54d61389-3d3a-412b-92c7-2eccc48bc909.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/dcc/dcc158e7-2ff6-4fa6-b51b-8d641af9-180b-42de-aaf0.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/a84/a84f2892-b177-4b65-b63a-4f1bdc6bd5a8.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/e74/e748067a-f7f5-4741-9eec-8748b45756ca.tiff |
| 308666239 | Mount Sinai Health System | 113 | NY | 69 | Not Hispanic or Latino | MALE | Unknown | Unknown | Chinese | Unknown | /sc/arion/projects/data-ark/digital_pathology_slides/e9c/e9c21f63-7bca-42bf-a9a3-22169e8e6c0e.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/051/05191f89-fcf8-43a5-8126-4b65a8b328.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/013/01300fde-0f8d-4f95-81c9-4651f8f8a0ba.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/220/2202fd96f-a586-4a33-9ae1-4633-9ae1-9f16387b0d69.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/c4d/c4d4a0a5-eef2-4d0b-9354-c21b704cf5.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/4b5/4b5dc7f0-7787-4260-919e-9b2a04acefa.tiff, /sc/arion/projects/data-ark/digital_pathology_slides/54d |

Pathology Data in AIR-MS



- Pathology case level metadata
- One entry per case

- Pathology report chunks as free text
- Full report is split into sections
 - Specimen Source
 - Grossing
 - Diagnosis
 - Comment
 - Addendum
- One entry per section

- Glass slide level metadata
- One entry per glass slide

- Digital slide level metadata
- One entry per digital slide

AIRMS Schema:
CDMPATHOLOGY

Case Metadata

Frequently used fields:

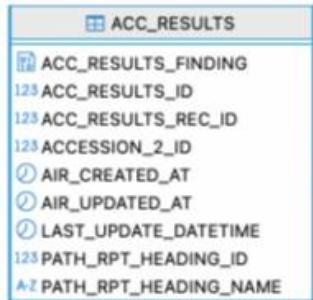
- *ACCESSION_2_ID* Linking key to other tables.
- *ACCESSION_NO* Accession number, used to identify cases. Pathologist-friendly.
- *MED_REC_NO* EPIC MRN. Can be used to link to other data modalities.

| ACCESSION | |
|-----------|---------------------------------|
| A-Z | ACC_CATG |
| ② | ACC_PROCESS_STEP_COMPLETED_DATE |
| 123 | ACCESSION_2_ID |
| A-Z | ACCESSION_NO |
| ② | AIR_CREATED_AT |
| ② | AIR_UPDATED_AT |
| ② | BIRTH_DATE |
| ② | CREATED_DATE |
| 123 | CURRENT_STATUS_ID |
| A-Z | FACILITY_CODE |
| 123 | FACILITY_ID |
| A-Z | FACILITY_NAME |
| A-Z | IMPORTED_CASE |
| ② | LAST_UPDATE_DATETIME |
| A-Z | MED_REC_NO |
| A-Z | MRN_FACILITY_CODE |
| A-Z | MRN_FACILITY_DESCRIPTION |
| A-Z | ORDER_NUMBER |
| A-Z | PATIENT_AGE |
| 123 | PATIENT_ID |
| A-Z | PERSONNEL_2_FULL_NAME |
| 123 | PERSONNEL_2_ID |
| A-Z | PROCESS_STEP_DESCRIPTION |
| A-Z | VISIT_NUMBER |

Report Text

Frequently used fields:

- *ACCESSION_2_ID* Linking key to case table.
- *PATH_RPT_HEADING_NAME* Name of the section. Section naming is not standardized. Some examples:
 - DIAGNOSIS
 - FINAL DIAGNOSIS
 - SPECIMEN SOURCE
 - NOTE
 - ADDENDUM 1, 2, 3, ...
 - COMMENT
 - GROSSING
- *ACC_RESULT_FINDING* Actual report text.



| | ACC_RESULTS |
|-----|-----------------------|
| | ACC_RESULTS_FINDING |
| 123 | ACC_RESULTS_ID |
| 123 | ACC_RESULTS_REC_ID |
| 123 | ACCESSION_2_ID |
| 🕒 | AIR_CREATED_AT |
| 🕒 | AIR_UPDATED_AT |
| 🕒 | LAST_UPDATE_DATETIME |
| 123 | PATH_RPT_HEADING_ID |
| A-Z | PATH_RPT_HEADING_NAME |

Glass Slide Metadata

| ACC_SLIDE | |
|-------------------------------------|---------------------------------|
| 123 | ACC_BLOCK_ID |
| A-Z | ACC_BLOCK_LABEL |
| ② | ACC_PROCESS_STEP_COMPLETED_DATE |
| 123 | ACC_SLIDE_ID |
| A-Z | ACC_SPECIMEN_DESCRIPTION |
| 123 | ACC_SPECIMEN_ID |
| 123 | ACCESSION_2_ID |
| ② | AIR_CREATED_AT |
| ② | AIR_UPDATED_AT |
| <input checked="" type="checkbox"/> | BIOPSY |
| ② | COLLECTION_DATE |
| A-Z | CONSULT_LABEL |
| A-Z | LAB_PROCEDURE_CODE |
| A-Z | LAB_PROCEDURE_DESCRIPTION |
| 123 | LAB_PROCEDURE_ID |
| ② | LAST_UPDATE_DATETIME |
| ② | RECV_DATE |
| A-Z | SLIDE_LABEL |
| 123 | SLIDE_NO |
| A-Z | SLIDE_TYPE |
| A-Z | SOURCE_MATERIAL_LABEL |
| A-Z | SOURCE_REC_TYPE |
| 123 | SPECIMEN_CATEGORY_ID |
| A-Z | SPECIMEN_CATEGORY_NAME |
| A-Z | SPECIMEN_GROUPS_CODE |
| 123 | SPECIMEN_GROUPS_ID |
| A-Z | SPECIMEN_LABEL |
| A-Z | TYPE |

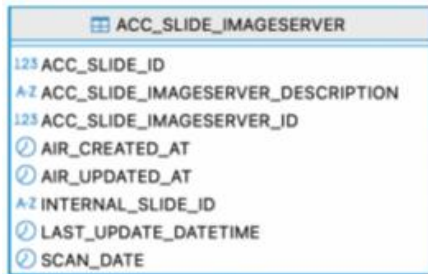
Frequently used fields:

- *ACCESSION_2_ID* Linking key to case table.
- *ACC_SLIDE_ID* Linking key to other tables.
- *SPECIMEN_LABEL* Specimen identifier within the case. Uppercase letter (A, B, ...). Matches specimen labels in free text pathology report. Pathologist-friendly.
- *ACC_BLOCK_LABEL* Block identifier within the specimen. Number (1, 2, ...). Matches block labels in free text pathology report, usually in the grossing description. Pathologist-friendly.
- *SLIDE_LABEL* Slide identifier within the block. Number (1, 2, ...). Not used in free text reports usually. Pathologist-friendly.
- *RECV_DATE* Received date of the specimen. Useful to match cases to other data modalities.
- *LAB_PROCEDURE_DESCRIPTION* Lab order that generated the slide. Usually a stain name. For example: H&E, PDL1, ER, PR, ...

Digital Slide Metadata

Frequently used fields:

- *ACC_SLIDE_ID* Linking key to slide table.
- *INTERNAL_SLIDE_ID* 6-char alphanumeric code that uniquely identifies a glass slide. Not pathologist-friendly. Aka BARCODE. If a barcode exists, the slide was scanned.
- *ACC_SLIDE_IMAGESERVER_DESCRIPTION* long alphanumeric code that uniquely identifies a .tiff file. Unique to a scanning event, not to a glass slide. Same barcode can be associated to multiple .tiff files (e.g., rescanning). Not pathologist-friendly.
- *SCAN_DATE* Timestamp of scanning event.



A screenshot of a database table schema for the table ACC_SLIDE_IMAGESERVER. The table has the following fields: ACC_SLIDE_ID (125), ACC_SLIDE_IMAGESERVER_DESCRIPTION (A-Z), ACC_SLIDE_IMAGESERVER_ID (125), AIR_CREATED_AT (timestamp), AIR_UPDATED_AT (timestamp), INTERNAL_SLIDE_ID (A-Z), LAST_UPDATE_DATETIME (timestamp), and SCAN_DATE (timestamp). Each field is preceded by a small circular icon containing a question mark.

| ACC_SLIDE_IMAGESERVER | |
|-----------------------|-----------------------------------|
| 125 | ACC_SLIDE_ID |
| A-Z | ACC_SLIDE_IMAGESERVER_DESCRIPTION |
| 125 | ACC_SLIDE_IMAGESERVER_ID |
| ⌚ | AIR_CREATED_AT |
| ⌚ | AIR_UPDATED_AT |
| A-Z | INTERNAL_SLIDE_ID |
| ⌚ | LAST_UPDATE_DATETIME |
| ⌚ | SCAN_DATE |

NOTE:

- Scanning started in 2020 at 40% capacity. 100% prospective digitization was reached in 2024. Most slides are still not scanned.
- *INTERNAL_SLIDE_ID* and *ACC_SLIDE_IMAGESERVER_DESCRIPTION* can be used to identify slides in Data Ark. Mapping is not publicly available.

Pathology Metadata – Practical Tips

- **How to communicate with the pathologist**

Accession numbers, specimen/block/slide labels are well understood by pathologists.

For example: slide “A1-1” from block “A1” in case MS-24-00001.

- **Linking slides to a diagnosis is challenging**

Reports are written for the entire case. Most granular diagnosis is at the specimen level (unstructured).

There is no slide level description. It is important to talk to a pathologist to understand which slides from a case/specimen are relevant to your study.

- **Some metadata can be inaccurate**

Some metadata fields exposed are used internally in the pathology department for workflow management. They are not for diagnostic purposes and can be inaccurate. Examples include: division, organ, specimen description, ICD codes. Always talk to a pathologist to discuss how to best approach your problem.

- **Data Ark slide links are not publicly available**

To obtain Minerva paths to a set of slides can be done through MSDW (requires IRB). You will need *INTERNAL_SLIDE_ID* and *ACC_SLIDE_IMAGESERVER_DESCRIPTION* for each slide in your cohort.

Pathologist collaboration is essential!

- Choosing a single representative image is challenging
 - Generally, requires subspecialty expertise by a pathologist
 - Context-dependent: a selected image for one project is not necessarily the best image for another project
- They can also make sure that the findings reported in the study make sense in the context of pathology
- For assistance in finding a pathologist interested in research, email the director of the Biorepository and Pathology CORE, Rachel Brody, MD, PhD (rachel.brody@mountsinai.org)

Summary

Summary & Conclusions

- Find more resources on the Minerva site:
<https://labs.icaahn.mssm.edu/minervalab/>
- We are in the process of developing more comprehensive links to the digital pathology files
- The slide viewer is under development
- If you have questions, you can submit a ticket here: <https://hpims.atlassian.net/servicedesk/customer/portal/67>