

# OMOP Abstractor

Michael Gurley

[m-gurley@northwestern.edu](mailto:m-gurley@northwestern.edu)

Northwestern Applied Research Informatics Group (NUARIG)

NUAIRG:

Firas Wehbe, MD, PhD

Yulia Bushmanova

NMEDW:

Daniel Schneider

Prasanth Nannapaneni

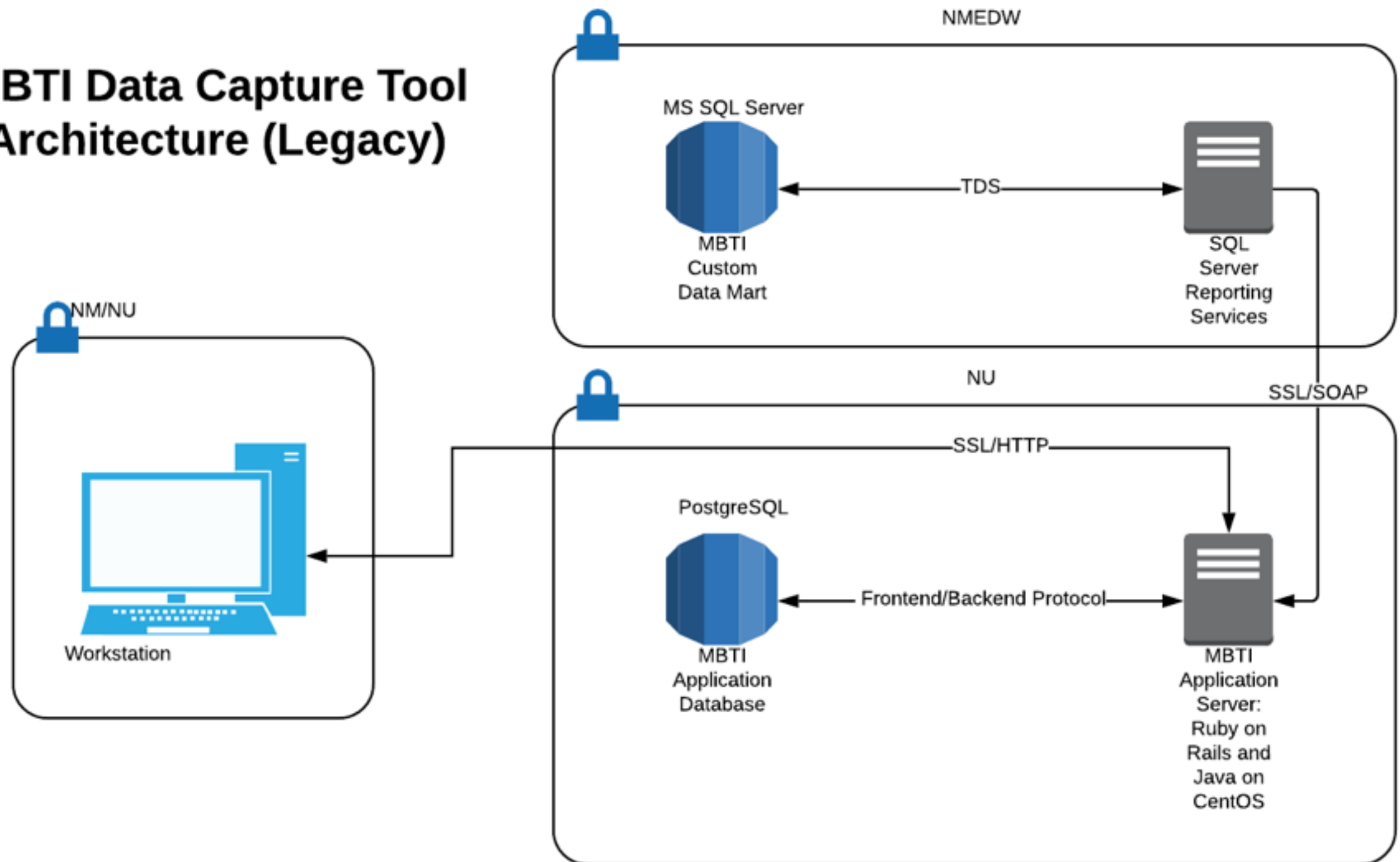
Martin Borsje

The Northwestern University logo, featuring the word "Northwestern" in a purple serif font, set against a white rectangular background.

# Background: MBTI Data Capture Tool

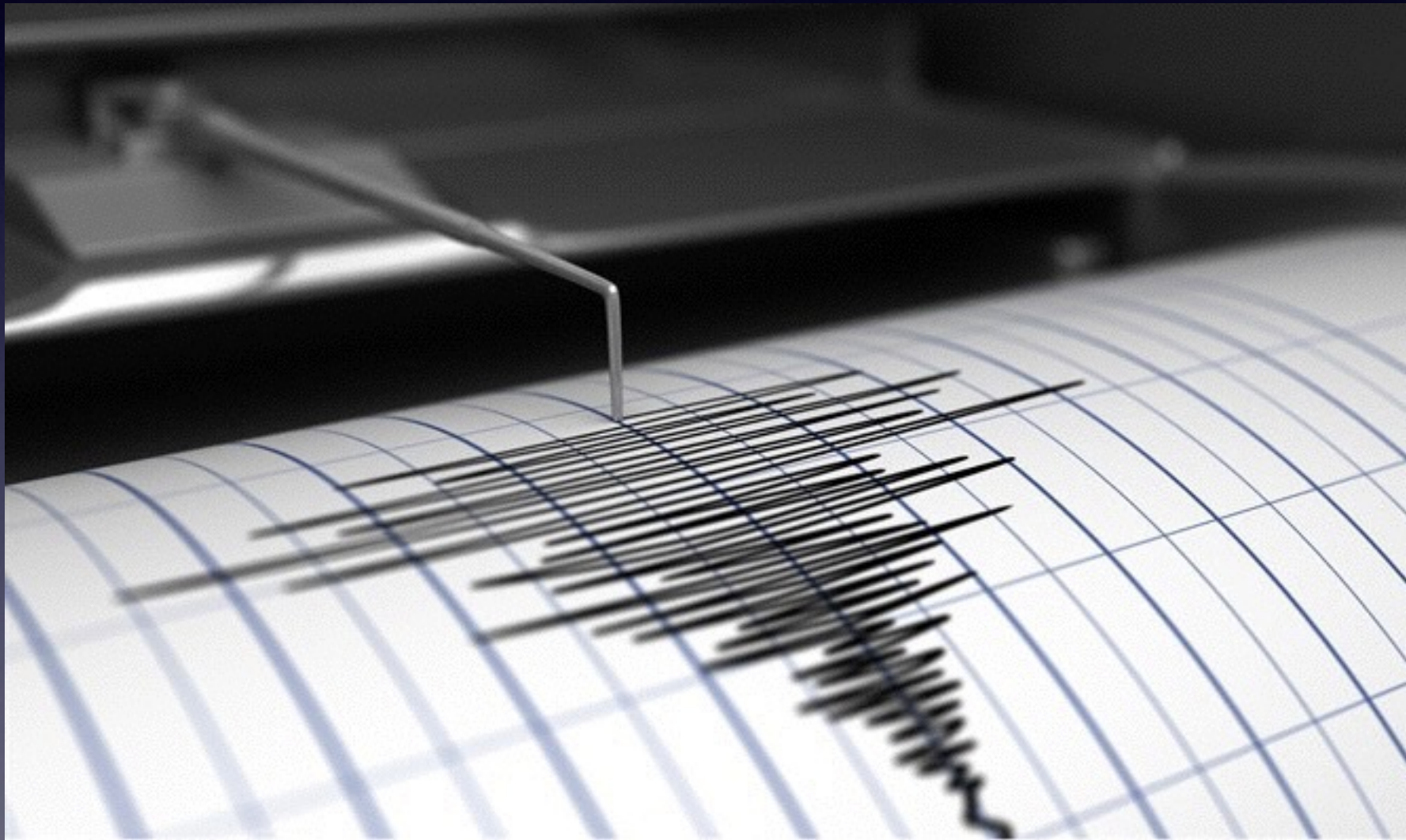
- The Northwestern Memorial Brain Tumor Institute outcomes/research database: The MBTI Data Capture Tool (MBTI-DCT).
- Backed by a custom ETL/data mart within the Northwestern Enterprise Data Warehouse (NMEDW) (Microsoft SQL Server and SQL Server Integration Services).
- Aggregates data from Northwestern Medicine's (NM) clinical systems: Epic (Outpatient), Cerner (Inpatient, Surgery, Pathology, Radiology) and MOSAIQ (Radiation Oncology).
- Custom data model.
- Incremental loads of the custom ETL/data mart into a PostgreSQL database via SQL Server Reporting Services (SSRS) exposed as SOAP Web Services.
- The MBTI-DCT includes an NLP-aided chart abstraction user interface for the curation of non-discrete data points from clinical narratives:
  - Anatomic site/histology, WHO grade, recurrence status, pathology findings from pathology reports.
  - Extent of Resection from surgical procedure reports.
  - Anatomical target of radiotherapy from radiation oncology summaries.
  - Performance status declarations, outside treatments from clinic progress notes.
  - Recurrence/progression declarations from imaging exam reports.
- Ruby on Rails user interface. Custom NLP pipeline using the Stanford NLP Java Library and Lingscope.

## MBTI Data Capture Tool Architecture (Legacy)





# Data Earthquake



# Change Happens

- NM merged with other hospital(s).
- Large project to migrate three Epic instances into one instance.
- NM migrated from using Cerner Surginet for tracking surgeries to Epic.
- NMEDW deprecated legacy traditionally modeled “Integrated Data Structures” to new fact/dimension modeled “Integrated Data Structures”.
- Programatic access to SSRS reports via SOAP services was disabled.

# Remediation Choice

- All this change broke the MBTI-DCT.
- The custom ETL/data mart's flow of new data stopped because of reliance on deprecated structures and data access strategies.
- Decision:
  - Remediate the ETL's population of the custom data model.

OR

- Replace the custom data model with a common data model (CDM) and remediate the MBTI-DCT UI and application logic to work with a CDM.
- All of Us and eMerge CDM activities were ramping up around this time. Laying the ground work for NM investing in transforming its NMEDW into the OMOP CDM.



# Choose OMOP: Challenges

- We decided to remediate the MBTI-DCT UI and application logic to work with the OMOP CDM.
- Challenges:
  - Work with a truncate/reload data refresh model. Our abstraction/NLP output/curation tables need stable structures to hang off of across data reloads.
  - Remediating MBTI-DCT UI to work with the OMOP CDM.
  - Representing PHI.
  - Include all surgeries in the OMOP instance.
  - Include all pathology procedures in the OMOP instance.
  - Include all pathology reports sections in the OMOP instance.
  - Preserve and represent references between surgeries, pathology procedures and pathology report sections.
  - De-duping surgeries and pathology procedures.

# Abstractor Library

- The legacy MBTI-DCT UI was backed by an NLP-aided assisted chart abstraction library, 'abstractor', created by NUAIRG.
- Abstractor's existing data model had wide support for specifying NLP expectations for chart abstraction. Beyond what is supported via the OMOP 'note\_nlp' table.
- Abstractor allows for the setup of 'expectations' necessary for NLP-aided assisted chart abstraction.
  - Event Cohort Expectations
    - Flows of clinical events. An abstractor 'namespace' is a set of saved search criteria to produce an 'event cohort'. For example: a 'Surgical Pathology' namespace that binds to all 'Final Diagnosis Section' notes related to 'Surgical pathology procedure' procedure occurrences
  - Data Point Expectations
    - Sets of abstractable data points for a specific type of clinical note or clinical event (subject). For example, ('Site', 'Histology', 'WHO Grade', 'IDH1 Status', etc.) for a clinical note (subject).
- Abstractor ties these two sets of expectations together to allow for an NLP library to generate suggestions for a defined set data points for a defined set of clinical events.
- Abstractor presents these flows of expectations (clinical events, data points and NLP suggestions) within a user interface for curation.



# Abstractor: Namespaces

- An abstractor namespace is a named search criteria that is instantiated within the 'abstractor\_namespaces' table. The search criteria are specified in the 'joins\_clause' and 'where\_clause' columns. The 'subject\_type' column indicates the subject or clinical event that is the target of the namespace. The 'id' of the 'abstractor\_namespaces' table can be referenced via in the 'namespace\_id' column of the 'abstractor\_subjects' table to bind data points to the 'namespace'.
- Abstractor has a schedulable 'suggestor' job that can cycle through all setup namespaces, execute each namespace's saved search criteria to return a list of subjects, discover all the data points points bound to the subject via the namespace and ask each data point's NLP configuration to generate NLP suggestions for the subject in hand.
- Bookkeeping that a subject for a given namespace has been processed is tracked in the 'abstractor\_namespace\_events' table.

# Abstractor: Subject, Predicate, Object

- Abstractor uses the familiar paradigm of subject, predicate, object to model data point expectations.
  - **Subject:** the 'abstractor\_subjects' table points to type of entity via the 'subject\_type' column and an abstractor namespace via the 'namespace\_id' column. For OMOP, normally the note table (or more precisely, as we will later see, the note\_stable\_identifier table) is the value of the 'subject\_type' column. The 'abstractor\_subjects' table also points to a 'predicate' via the 'abstractor\_abstraction\_schema\_id' column.
  - **Predicate:** the abstractor\_abstraction\_schemas table describes a property of the subject. For example, the 'has\_icdo3\_histology' property of declaring the histology of a cancer diagnosis. A predicate can have different types of possible answers. This is captured in the abstractor\_object\_type\_id column: list, number, boolean, string, date and text. Some predicates can have multiple syntactical expressions or synonyms. This is specified in the 'abstractor\_abstraction\_schema\_predicate\_variants' table.
  - **Object:** the 'abstractor\_object\_values' table describes, for 'list' predicates, the range of possible values that can satisfy the property of the subject. Some objects can have multiple syntactical expressions or synonyms. This is specified in the 'abstractor\_object\_value\_variants' table.
- For those data points that are only meaningfully collected in groups (like a cancer diagnosis 'site' and 'histology') and with possible cardinality > 1, abstractor allows for repeating grouping of data points via the 'abstractor\_subject\_groups' table.



# Abstractor: Sources

- Abstractor has the concept of 'sources' for each data point for each subject. The 'abstractor\_abstraction\_sources' table specifies in the 'from\_method' column the textual content of the subject. Normally, for OMOP this will be the 'note\_text' column of the 'note' table but could also be the 'note\_title' column.
- Each source entry must specify the expected format of how a predicate appears in content for the subject in the 'abstractor\_rule\_type\_id' column. Here are the possible rule types:
  - 'name/value': both the predicate and the value are expected within a suggestible mention. For example, 'KPS: 20' or 'The patient has a KPS of 20.'
  - 'value': only the value is expected within a suggestible mention. For example: 'Pituitary adenoma (see Microscopic Description).'
- Each source entry must specify a strategy for making suggestions in the 'abstractor\_abstraction\_source\_type\_id' column. Here are the possible strategies:
  - 'nlp suggestion': delegate the generation of suggestions to abstractor's built-in NLP library. A simple sentence-splitting, rule-based named entity recognizer with negation detection.
  - 'custom nlp suggestion': delegate the generation of suggestions to an external NLP library via communication over a RESTful interface. Not confined to rule-based NLP. Abstractor sends a document to the configured endpoint with an embedded endpoint URL containing the dictionary for each data point and a call back endpoint to send back suggestions.

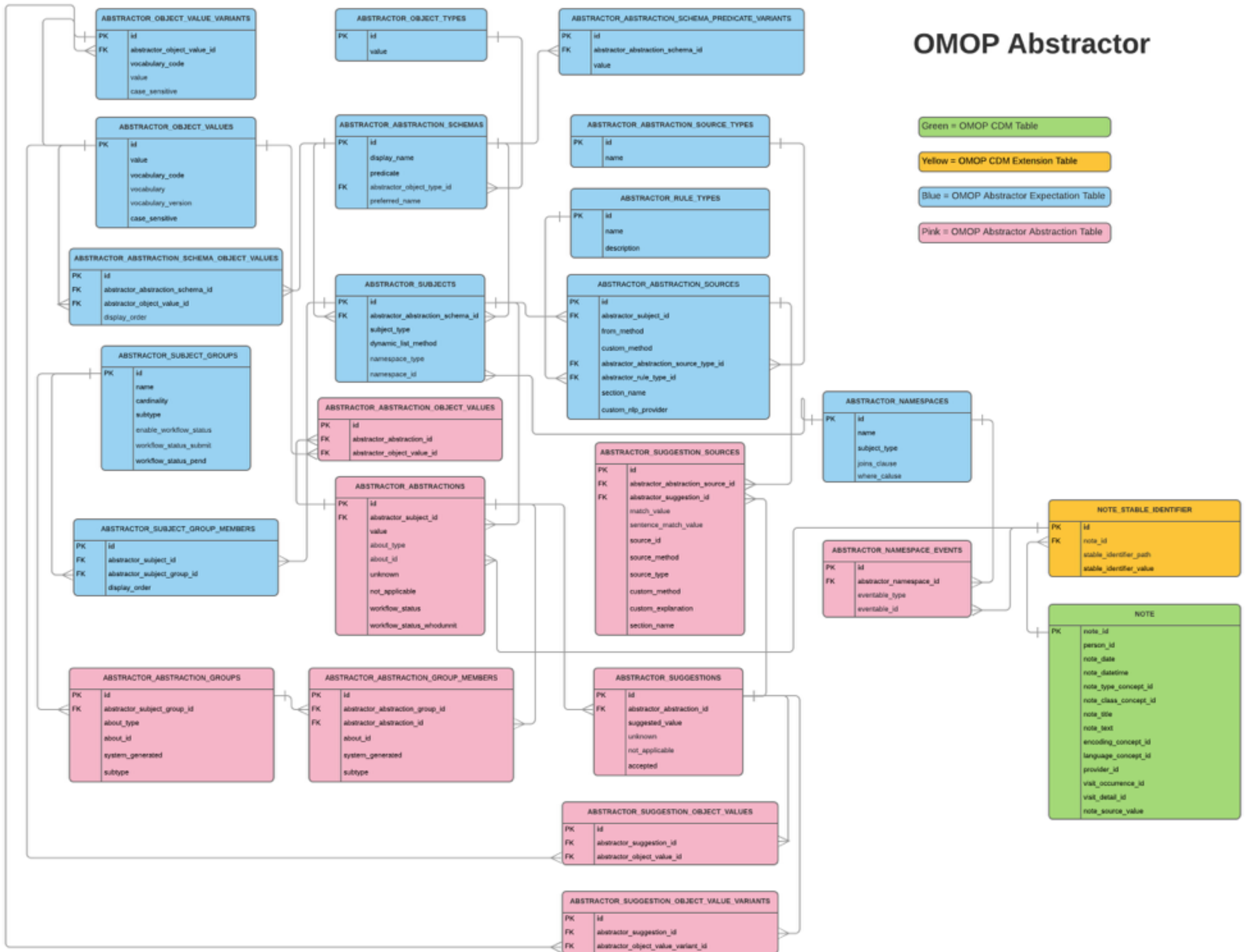


# Abstractor:

## Suggestions and Abstractions

- Abstractor has the concept of 'suggestions' for each data point for each subject. The 'abstractor\_suggestions' table contains the NLP suggested value or suggestion of 'unknown' or 'not applicable'. The underlying textual evidence supporting an NLP library's suggestion is contained within the 'match\_value' and 'sentence\_match\_value' columns of the 'abstractor\_abstraction\_sources' table. Abstractor automatically makes a suggestion of 'unknown' and 'not applicable' for all data points.
- A curator's accepted 'suggestion' for a data point is recorded within the 'abstractor\_abstractions' table. This table is the center of the abstractor universe.
- Grouped, relating accepted data points are collected within the 'abstractor\_abstraction\_groups' table.

## OMOP Abstractor



# Challenge:

## Work with a truncate/reload data refresh model.

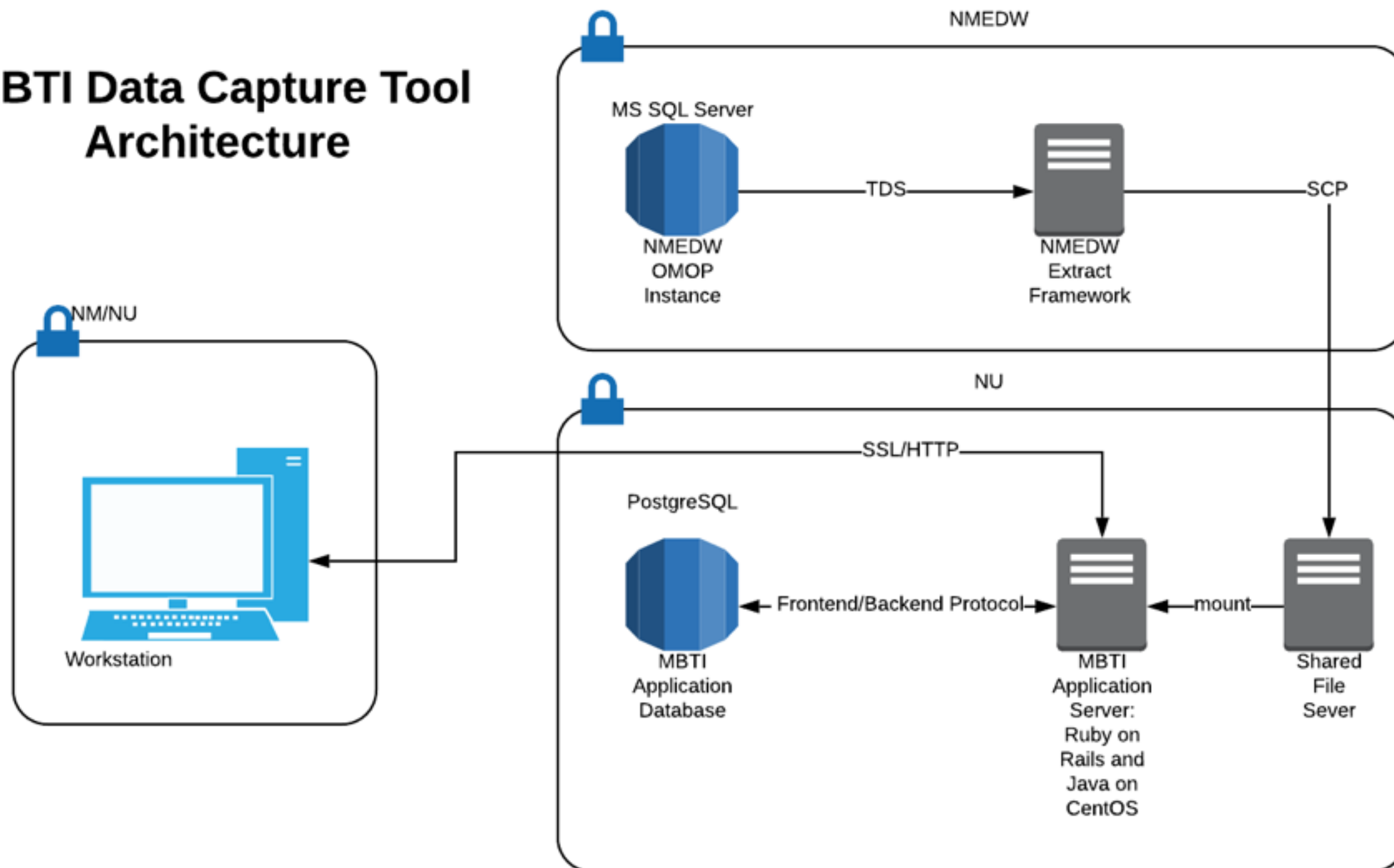
- New load strategy:
  - Replacement of programmatic access to data via SSRS.
  - NMEDW has a single OMOP instance that is partitioned during extraction by a cohort definition for the MBTI.
  - The NMEDW extract framework deposits files on a shared folder mounted to the MBTI-DCT application server.
  - Incremental loads not feasible.
- Stable Identifiers
  - OMOP internal IDs change across truncate/reload data refreshes.
  - Our abstraction/NLP output/curation tables need stable structures to hang off of across data reloads.
  - Asked the NEMEDW OMOP data architect team to populate stable identifier tables for NOTE and PROCEDURE\_OCCURRENCE.
  - The stable identifier tables contains an invariant 'id' column that stays stable across loads, an OMOP internal id column ('note\_id', 'procedure\_occurrence\_id') that changes across loads and invariant pointers to source-system row-level provenance via the 'stable\_identifier\_path' and 'stable\_identifier\_value' columns.



# NOTE\_STABLE\_IDENTIFIER

```
CREATE TABLE public.note_stable_identifier  
(  
  id bigint NOT NULL DEFAULT nextval('note_stable_identifier_id_seq'::regclass),  
  note_id bigint NOT NULL,  
  stable_identifier_path character varying NOT NULL,  
  stable_identifier_value character varying NOT NULL,  
  CONSTRAINT note_stable_identifier_pkey PRIMARY KEY (id)  
)
```

# MBTI Data Capture Tool Architecture



# Challenge:

## Remediating MBTI-DCT UI to work with the OMOP CDM

### Index Page

- Legacy MBTI-DCT UI was based on a custom data model. Each type of note targeted for NLP-aided chart abstraction had its own dedicated table: pathology reports, imaging exam reports, clinic progress notes and radiation oncology summaries. Each note type had its own dedicated index and edit page.
- Index Page
  - Remediate one index screen to display notes needing curation by abstraction 'namespaces'.
  - Remediate the index screen to allow for searching:
    - By keyword search across OMOP PHI tables and the note text.
    - By note date.
    - By providers associated with a first level procedures (for example, pathology procedures) and second level procedures (for example, surgeries).



Challenge:  
Remediating MBTI-DCT UI  
to work with the OMOP CDM  
Edit Page

- Edit Page
  - Remediate the edit screen to display data from the OMOP PHI tables.
  - Remediate the edit screen to display the list of abstractable data points for a note based on the set of abstractable data points bound to a 'namespace'.
  - Remediate the edit screen to display associated first level and second level procedures.
  - Remediate the edit screen to display other note entries associated to first level procedures (for example, all other "sections" of the current pathology report section)

# Northwestern MBTI Data Capture Tool

## Notes

Reviewed?  
needs review

Provider  
x HORBINSKI, CRAIG

Namespace  
Molecular Pathology  
Surgical Pathology  
Outside Surgical Pathology

Search

From

To

SEARCH

Clear

< 1 2 3 4 5 6 7 8 9 10 11 12 >

Note Date	Note Type	Note Title	First Name	Last Name	MRN(s)	
01/01/1900	Note	Final Diagnosis	Bob	Jones	Northwestern 000000000	Review
01/01/1900	Note	Final Diagnosis	Bob	Jones	Northwestern 000000000	Review
01/01/1900	Note	Final Diagnosis	Bob	Jones	Northwestern 000000000	Review
01/01/1900	Note	Final Diagnosis	Bob	Jones	Northwestern 000000000	Review

Northwestern MBTI Data Capture Tool

# Note

[Back](#) | [Notes](#) | [Previous](#) | [Next](#)

Patient	MRN(s):	Note Date	Note Type	Note Class	Title	Provider
Bob N Jones	Northwestern 000000000	01/01/1900	Note	No matching concept	Final Diagnosis	

## Procedures

Procedure	Date	Provider	Specimens	Notes
Surgical pathology procedure	01/01/1900	HORBINSKI, CRAIG M.		Specimen/Gross Description <div>VIEW</div>
				Intraoperative Consultation Findings <div>VIEW</div>
				Clinical Information <div>VIEW</div>
				Surg Path Non-Chartable Comment <div>VIEW</div>
				Addendum <div>VIEW</div>
Craniectomy, trephination, bone flap craniotomy; for excision of meningioma, supratentorial	01/01/1900	CHANDLER, JAMES P.		



NOT APPLICABLE ALL

UNKNOWN ALL

## Metastatic Cancer

### Histology

- ☒ adenocarcinoma, metastatic (8140/6) 
- ☐ not applicable
- ☐ unknown

Edit | CLEAR

### Site

- ☐ cerebellum, nos (c71.6) 
- ☐ not applicable
- ☐ unknown

Edit | CLEAR

### Primary Site

- ☐ not applicable
- ☐ unknown

Edit | CLEAR

### Laterality

- ☐ not applicable
- ☐ unknown

Edit | CLEAR

### Recurrent

- ☐ not applicable
- ☐ unknown

Edit | CLEAR

NOT APPLICABLE GROUP

UNKNOWN GROUP

ADD METASTATIC CANCER

## Note text

A and B. Tumor, cerebellum, resection:

**Metastatic adenocarcinoma** (see Note).

**Note:** This tumor shows tall columnar cells with luminal necrosis, the combination of which is a classic hallmark of colorectal adenocarcinoma.

IDH1 Status

☐ not applicable

☐ unknown

Edit

CLEAR

IDH2 Status

☐ not applicable

☐ unknown

Edit

CLEAR

1P Status

☐ deleted

☐ non-deleted

SAVE

Cancel

19q Status

☐ not applicable

☐ unknown

Edit

CLEAR

10q/PTEN Status

☐ not applicable

☐ unknown

Edit

CLEAR

MGMT promoter methylation status Status

☐ not applicable

☐ unknown

Edit

CLEAR

ki67

☐ not applicable

☐ unknown

Edit

CLEAR

p53

☐ not applicable

☐ unknown

Note text

A and B. Tumor, cerebellum, resection:  
Metastatic adenocarcinoma (see Note).

Note: This tumor shows tall columnar cells with luminal necrosis, the combination of which is a classic hallmark of colorectal adenocarcinoma.

# Challenge:

## PHI

- The MBTI-DCT displays PHI. The display of PHI within the MBTI-DCT is necessary to meet curation and outcomes use cases.
- Solution: adopted and reused the PHI table specified and populated for the All of Us project.
  - pii\_address
  - pii\_email
  - pii\_mrn
  - pii\_name
  - pii\_phone\_number



# Challenge: Include all surgeries

- Validate that surgeries were **not** being included in the current OMOP build.
- Surgeries at NM are tracked in the new consolidated Epic instance.
- For Epic, the current OMOP build was populating the PROCEDURE\_OCCURRENCE table exclusively from charge-oriented tables
  - HSP\_ACCT\_PX\_LIST
  - HSP\_ACCT\_CPT\_CODES
  - HSP\_TRANSACTIONS
  - ARPB\_TRANSACTIONS
- Asked the NEMEDW OMOP data architect team to pull from actual Epic surgery tables.
  - or\_log
  - or\_log\_all\_proc
  - or\_proc
  - or\_proc\_cpt\_id
- Automatic mapping to Procedure domain standardized vocabulary entries was achieved by using the or\_proc\_cpt\_id.real\_cpt\_code field.

Opinion Sidebar:  
Prefer Small (accurate) Data  
over  
Big (messy) Data

- Simple Determinism is better than Clever Probabilism.
  - If a source system represents a class of clinical events in a discrete manner to support a clinical workflow (like Epic does for surgeries), prefer this canonical representation to the clinical event welter caused by charge-oriented representations.
- Open question:
  - Exclude charge-related representations?

**YEAH, WELL, YOU KNOW,**



**THAT'S JUST, LIKE,  
YOUR OPINION, MAN.**



# Challenge:

## Include all pathology procedures

- Validate that pathology procedures were **not** being included in the current OMOP build.
- Pathology procedures at NM are tracked in Cerner Pathnet Anatomic Pathology.
- For Cerner, the current OMOP build was populating the PROCEDURE\_OCCURRENCE table from charge-oriented tables
  - encounter (from Cerner)
  - EPSI charge tables
- Asked the NEMEDW OMOP data architect team to pull from actual Cerner Pathnet Anatomic Pathology tables.
  - pathology\_case
  - prefix\_group
  - case\_specimen
  - case\_report
  - clinical\_event
  - ce\_blob
- Made mappings from local prefix\_group entries to standardized Procedure domain entries in the SNOMED vocabulary. Mostly along this axis: Procedure | Laboratory Procedure (procedure) | Anatomic Pathology Procedure. See spreadsheet.
- Need to map local Cerner case\_specimen specimen codes to the Specimen domain entries in the SNOMED vocabulary.

# Challenge:

## Include all pathology reports sections in the OMOP instance.

- Validate that pathology reports sections were being included in the current OMOP build ***as separate entries in the NOTE table with the section name populating the note\_title field.***
- Pathology reports sections at NM are tracked in Cerner Pathnet Anatomic Pathology.
- Pathology reports are written in 'sections'. Each section having a dedicated purpose. For example: 'Final Diagnosis', 'Microscopic Description', 'Specimen/Gross Description' and 'Clinical Information'.
- Most often the data points desired to be extracted from a pathology report reside in the 'Final Diagnosis' section.
- Other 'sections' can often be the source of false positives for NLP pipelines. For example, historical diagnoses mentioned in the 'Clinical Information' section.
- Simple Determinism is better than Clever NLP.
  - If a source system splits a pathology reports into discrete labeled sections to support a clinical workflow (like Cerner does for pathology reports), prefer this canonical representation instead of a multi-section conglomerated representation.
- Don't use the sectionizing component of your NLP pipeline if your source system sectionizes for you.

# Challenge:

## Preserve and represent references between surgeries, pathology procedures and pathology report sections.

- Validate that references between surgeries, pathology procedures and pathology report sections were **NOT** being included in the current OMOP build.
- The conventional advice to tie OMOP clinical events by joining to VISIT\_OCCURRENCE **is insufficient**.
  - Possible for one VISIT\_OCCURRENCE to span multiple surgery entries and multiple pathology procedures in PROCEDURE\_OCCURRENCE.
  - Possible for one VISIT\_OCCURRENCE entry to span multiple pathology reports.
- Make explicit references between surgical PROCEDURE\_OCCURRENCE entires and pathology PROCEDURE\_OCCURRENCE entires via FACT\_RELATIONSHIP.
  - Asked the NMEDW OMOP data architect team to build a join table within the NMEDW integrated data structures associating Cerner pathology procedures and Epic surgeries. Match on patient and surgery date to pathology accession date/case collection date. ETL the join table into the OMOP FACT\_RELATIONSHIP table
- Make explicit references between pathology report section NOTE entires and pathology PROCEDURE\_OCCURRENCE entires via FACT\_RELATIONSHIP (we are not on OMOP CDM version yet where this can be done directly within the NOTE table with note\_event\_id and note\_event\_field\_concept\_id).
  - Asked the NMEDW OMOP data architect team to build a join table within the NMEDW integrated data structures associating Cerner pathology procedures and Cerner pathology report sections. ETL the join table into the OMOP FACT\_RELATIONSHIP table



```

SELECT note.note_id~
~~~~, note.note_date~
~~~~, note_stable_identifier.id~
~~~~, note_stable_identifier.stable_identifier_path~
~~~~, note_stable_identifier.stable_identifier_value~
~~~~, note.note_title~
~~~~, note.note_text~
~~~~, procedure_occurrence.procedure_occurrence_id~
~~~~, procedure_occurrence.procedure_concept_id~
~~~~, concept.concept_code~
~~~~, procedure_occurrence.procedure_date~
~~~~, procedure_occurrence_stable_identifier.id~
~~~~, procedure_occurrence_stable_identifier.stable_identifier_path~
~~~~, procedure_occurrence_stable_identifier.stable_identifier_value_1~
~~~~, prov1.provider_name~
~~~~, prov2.provider_name~
~~~~, pos12.id~
~~~~, pos12.stable_identifier_path~
~~~~, pos12.stable_identifier_value_1~
FROM note_stable_identifier JOIN note~
~~~~~ ON stable_identifier_value_1.note_id = note.note_id~
~~~~~ JOIN fact_relationship~
~~~~~ ON fact_relationship.domain_concept_id_1 = 5085 AND fact_relationship.fact_id_1 = note.note_id AND fact_relationship.relationship_concept_id = 44818790~
~~~~~ JOIN procedure_occurrence~
~~~~~ ON fact_relationship.domain_concept_id_2 = 10 AND fact_relationship.fact_id_2 = procedure_occurrence.procedure_occurrence_id AND procedure_occurrence.procedure_concept_id = 4213297~
~~~~~ JOIN procedure_occurrence_stable_identifier~
~~~~~ ON procedure_occurrence.procedure_occurrence_id = procedure_occurrence_stable_identifier.procedure_occurrence_id~
~~~~~ JOIN concept~
~~~~~ ON procedure_occurrence.procedure_concept_id = concept.concept_id~
~~~~~ JOIN fact_relationship AS fr2~
~~~~~ ON fr2.domain_concept_id_1 = 10 AND fr2.fact_id_1 = procedure_occurrence.procedure_occurrence_id AND fr2.relationship_concept_id = 44818888~
~~~~~ JOIN procedure_occurrence pr2~
~~~~~ ON fr2.domain_concept_id_2 = 10 AND fr2.fact_id_2 = pr2.procedure_occurrence_id~
~~~~~ JOIN procedure_occurrence_stable_identifier pos12~
~~~~~ ON pr2.procedure_occurrence_id = pos12.procedure_occurrence_id~
~~~~~ JOIN provider prov1~
~~~~~ ON procedure_occurrence.provider_id = prov1.provider_id~
~~~~~ JOIN provider prov2~
~~~~~ ON pr2.provider_id = prov2.provider_id~
WHERE note.note_title = 'Final Diagnosis'

```

# FACT\_RELATIONSHIP

- Entry between pathology report section in NOTE and pathology procedure in PROCEDURE\_OCCURRENCE:  
relationship\_concept\_id = 44818790  
'Has procedure context (SNOMED).' Plus converse entry.
- Entry between pathology procedure in PROCEDURE\_OCCURRENCE and surgery in PROCEDURE\_OCCURRENCE:  
relationship\_concept\_id = 44818888  
'Procedure context of (SNOMED)' Plus converse entry.

# Challenge:

## De-duping surgeries and pathology procedures

- Entries in FACT\_RELATIONSHIP allow us to surface these clinical events from the morass of charge-related representations for the same clinical events. Enabling de-duplication.
- Would be nice if OMOP natively contained some kind of way of designating entries in the PROCEDURE\_OCCURRENCE table as canonical or first-class versus entries that are financial echoes of the same clinical events



# Future

- Change and Challenges:
  - Extract and open source the mashup of 'abstractor' and 'OMOP' as "OMOP Abstractor". Coming Soon!
    - <https://github.com/NUARIG/omop-abstractor>
  - Integrate the OMOP vocabulary tables.
  - Need to incorporate into our OMOP instance legacy surgeries from Cerner Surginet.
  - Need to incorporate into our OMOP instance pathology procedures from a Cerner Co-Path instance to be loaded into our NMEDW.
  - Need to incorporate into our OMOP instance pathology procedures from a pending migration to Epic Beaker.
  - Improving our NLP algorithms. Not an NLP programmer. NLP pipeline has a RESTful interface that can delegate the generation of suggestions for a document and namespace to an endpoint and receive back suggestions via a endpoint. So better NLP can be used.

Thanks!