# Ontologies, Data Curation and Text Mining

Thérèse Vachon, Global Head of Text Mining Services

Novartis Institutes for Biomedical Research, NX

Information Retrieval and Text Mining for Biology
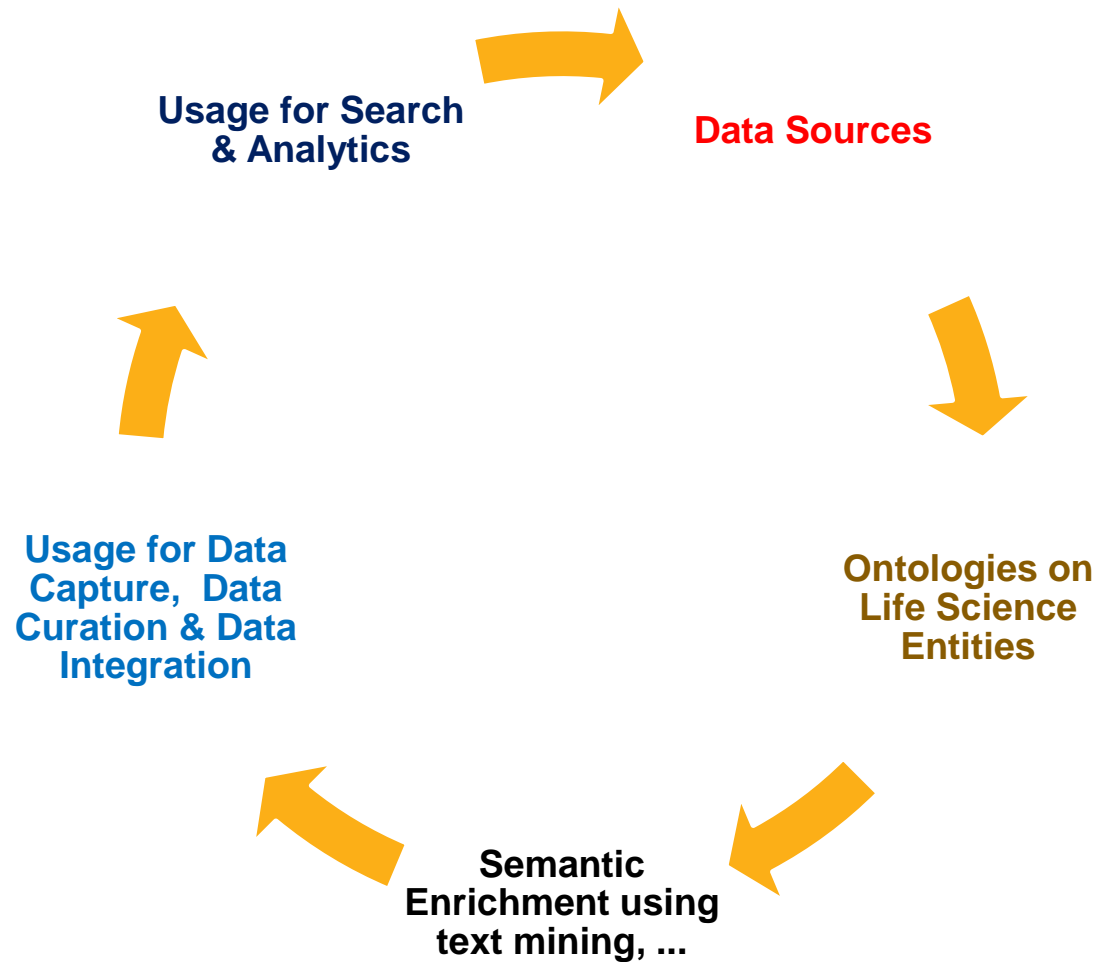
Geneva, 05.06.2015

**U NOVARTIS**

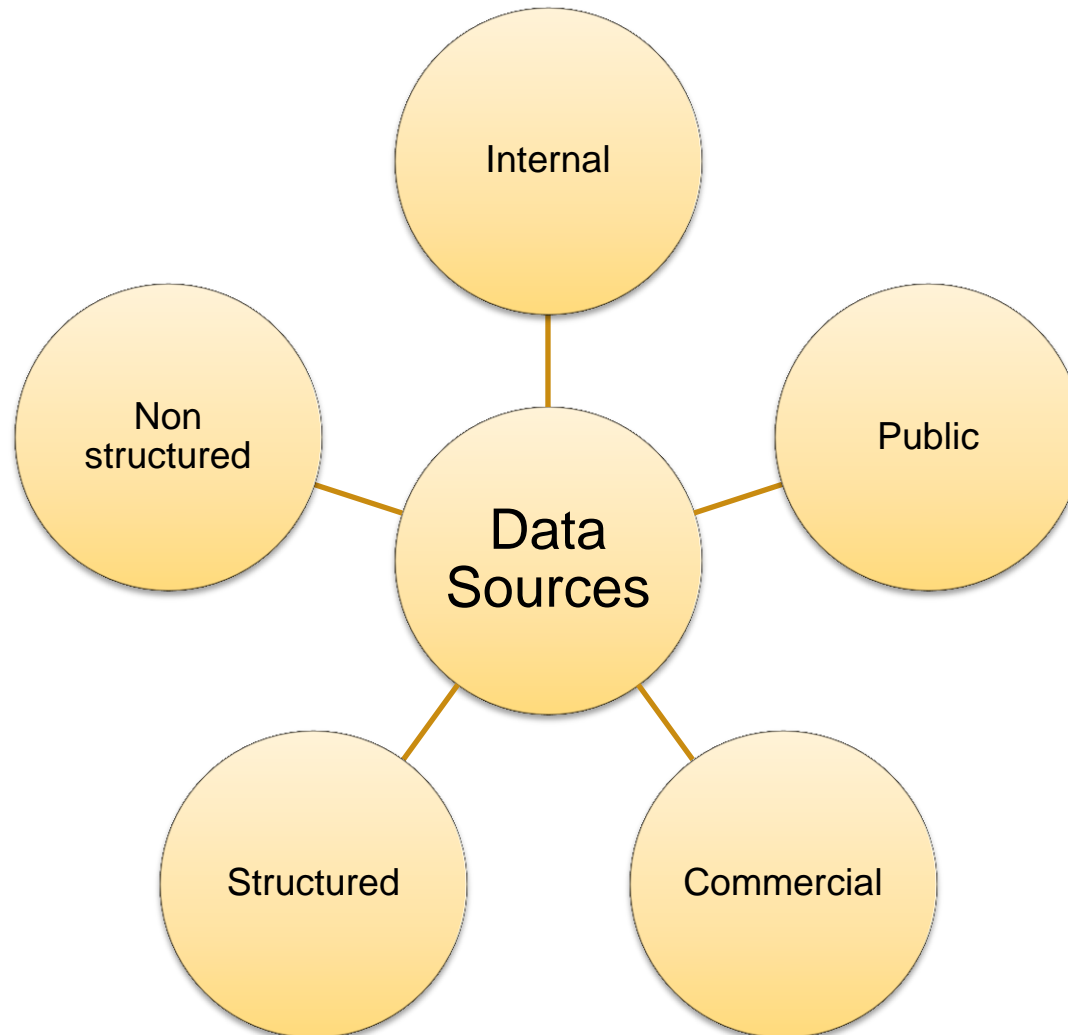# Main activities in Text Mining Services

- NIBR Ontologies

- Scientific Data Curation / Data Integrity & Consistency

- Scientific Data Integration

- RDF Graph DBs

- Federated Queries

- Search

- Text Mining

- Patent Mining

U NOVARTIS

# Agenda

Usage for Search & Analytics

Data Sources

Ontologies on Life Science Entities

Semantic Enrichment using text mining, ...

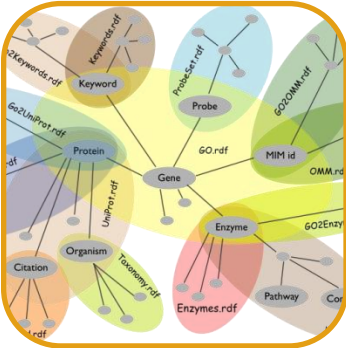Usage for Data Capture,  Data Curation & Data Integration
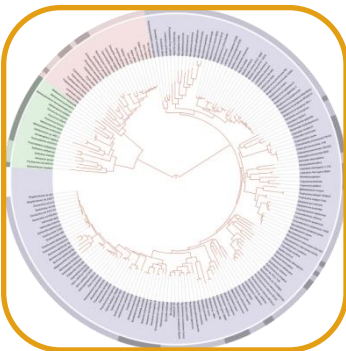
NOVARTIS

# Data Sources

# Central Repository for NIBR Ontologies



It aims at providing a uniform vocabulary within and across data repositories and ontology services to informatics teams for usage within NIBR applications.



Why is Ontology good for science?

- Standardized vocabulary with definitions and synonyms for unified database annotations across different databases
- Hierarchical organization for aggregation and multi-level comparison of results
- Community adoption for easy comparison of results to other project results worldwide
- Explicit relationships and underlying logical definitions for automated reasoning
- Explicit bridging relationships between different ontologies for exploring underlying mechanisms

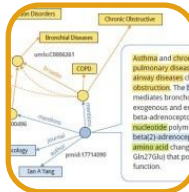# NIBR Ontologies, domain of applications

It can be used for multiple purposes e.g.


Registration systems (metadata capture), data curation


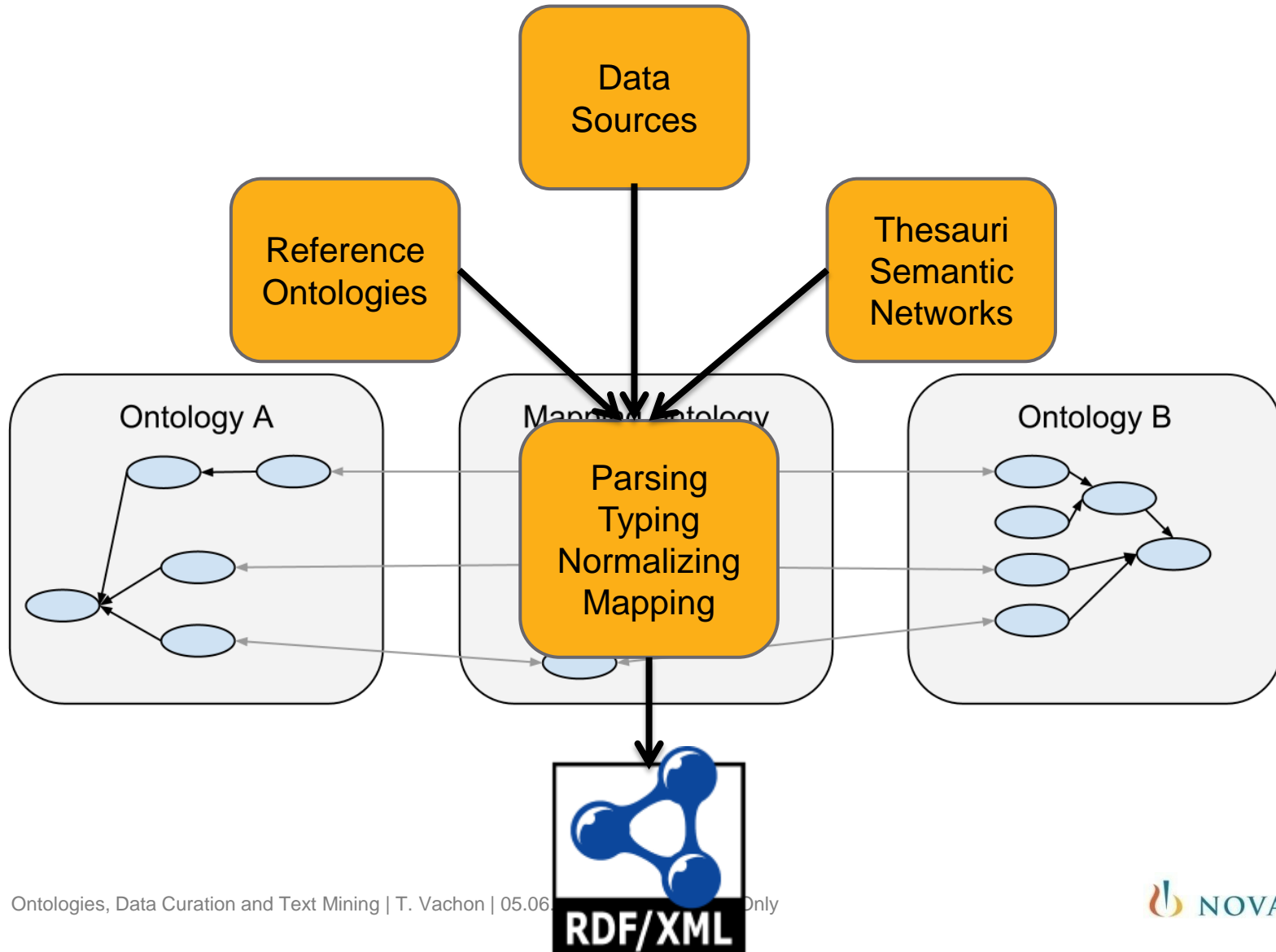Data mapping, bridging ontologies, navigation between concepts and referential data
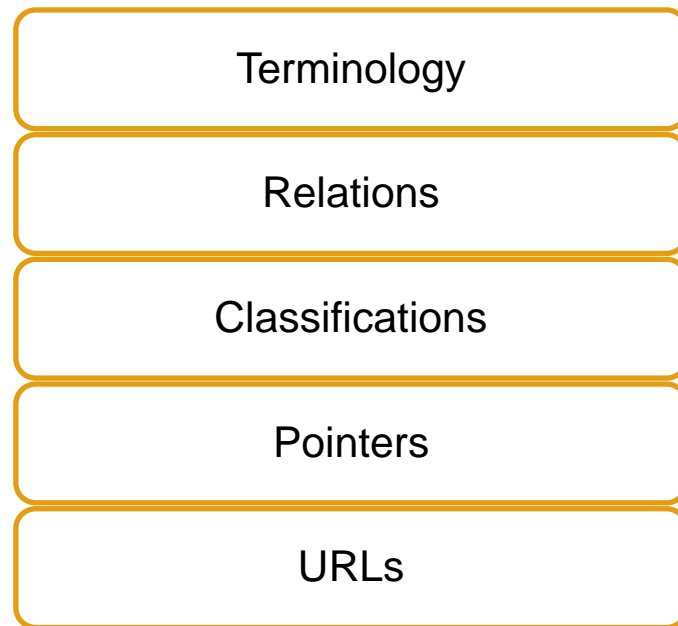

Text and data mining


Contribution to scientific queries, semantic data federation

NOVARTIS

# NIBR Ontologies, Content Generation

NOVARTIS

# NIBR Ontologies, Organization by concept type

Terminology

Relations

Classifications

Pointers

URLs

NOVARTIS

# NIBR Ontologies, Concept Types



Genes
Proteins
Pathways
Species
Target Classes

Cell lines
Cell types
Anatomy

Diseases

Compounds,
Products

Companies
Public
Institutions
Countries

ADCs
Instruments
Skills

NOVARTIS

# Auto-Suggest

# Bridging Ontologies

# Classifications

## Hierarchical Classifications
### Example on Diseases



(21/38211) DISEASES
- (8/354) abdominal disorder
- (14/1127) cardiovascular disorder
- (8/2969) digestive system disorder
- (5/201) ear, nose and throat disorder
- (43/1298) eye disorder
- (48/947) genetic disorder
- (11/1233) hematologic disorder
- (8/392) immune system disorder
- (11/668) infant or newborn disorder
- (2/289) injury and intoxication
- (2/3465) mental and neurological disorder
- (17/1061) musculoskeletal disorder
- (4/10856) neoplasm
- (2/2163) nutritional and metabolic disorder
- (112/5760) pathological conditions, signs and s
- (20/850) respiratory tract disorder
- (8/1418) skin and connective tissue disorder
- (2/670) therapy, prevention and control
- (2/36) thoracic disorder
- (17/24) Unclassified Terms
- (9/2409) urogenital tract disorder

Metastore Technology — a spectrum of options 4.4.2014, T. Vachon

(11/1233) hematologic disorder
- (10/160) blood clotting disorder
- (13/83) blood protein disorder
- (12/102) bone marrow disorder
  - (2/2) aplastic anemia
  - bone marrow aplasia
  - bone marrow depression
  - bone marrow neoplasm
  - bone marrow suppression
  - chemotherapy induced bone marrow injury
  - (1/2) hypoplastic anemia
  - (10/16) leukopenia
  - (8/8) myelodysplastic syndrome
  - myelodysplastic/myeloproliferative disease
  - (8/62) myeloproliferative disorder
  - pancytopenia
- (6/80) erythrocyte disorder
- hematologic genetic disorders
- hematologic pregnancy complications
- (4/191) hematological neoplasm
- (1/1) hematopoiesis disorder
- (2/2) hemophagocytic lymphohistiocytosis
- (9/222) leukocyte disorder
- (12/381) lymphatic system disorder

NOVARTIS

# Navigation

# Annotation using text mining and NIBR Ontologies

# Text extractors

- **Lexical extraction**

- **Pattern extraction**
  - E.g. IDs, Patent numbers, Compound numbers

- **Chemical Entity Extraction**
  - Trivial names
  - IUPAC
  - Smiles
  - Inchi
  - Images

- **Bulk Terminology Extractor (see examples)**

# More sophisticated text mining

- Extraction of entities relationships from literature

- Extraction of skills and disambiguation of authors from literature

- Text mining projects specific for a team

# Data Curation

# Scientific Data Curation
## *Curation types, Domains*

U NOVARTIS

# Expectations about scientific data curation in NIBR

- Implementing rules and best practices for data and metadata capture

- Helping to register / publish new data in the existing systems

- Adding / correcting metadata using authoritative vocabularies

- Identifying missing terms / terminologies

- Correcting faulty data in the source systems

- QA/QC: Assessing the data integrity in different systems & implementing corrective measures

- Making sure that the source data is consumable by the systems

- Curating different entities (assays, targets, compounds, genomics data, etc)

- Curating internal / external data sources (literature, patents, etc)

NOVARTIS

# Curation Framework

*What is the CTMF ?*

- **NIBR Ontologies (Metastore)** provide an integrated terminology and synonyms resolution service.

- The **C**ollaborative **T**erminology **M**anagement **F**ramework enables a distributed creation of the Metastore content.

- In its **first version**, the CTMF provides the following functionalities:

  - Users can request a new term to be added

  - The CTMF supports content owners and users in the clarification and resolution process.

  - Applications properly linked to the CTMF can make use of "temporary ID" before the term resolution process is completed.

# CTMF: an introduction
*Request for new terms*

- A user can request a new term via the **CTMF application**, or in any application that properly includes the **CTMF widget**.

- **The CTMF is not a registration systems**: no automatic bulk submissions of terms are possible, as each term is validated by content owners.

# CTMF: an introduction
## *Clarification and resolution process*



- The CTMF presents a "term status page" where discussion on the meaning of the term can be recorded.

- In general, content owners can "resolve" a new term as relative to a new concept, as a synonym of an existing concept, or as an error.

# CTMF: an introduction
## *Use of temporary IDs by applications*



Within a "CTMF enabled application", the user requests a new term (e.g. a new species in a registration system)

**Submitter**

"my new species"

**CTMF**

MVNTMPXXXX

The CTMF releases a "temporary ID", that can be used "consistently" by the application to refer to the new term (e.g.: a new species). The same temporary ID will be released for the same "term".
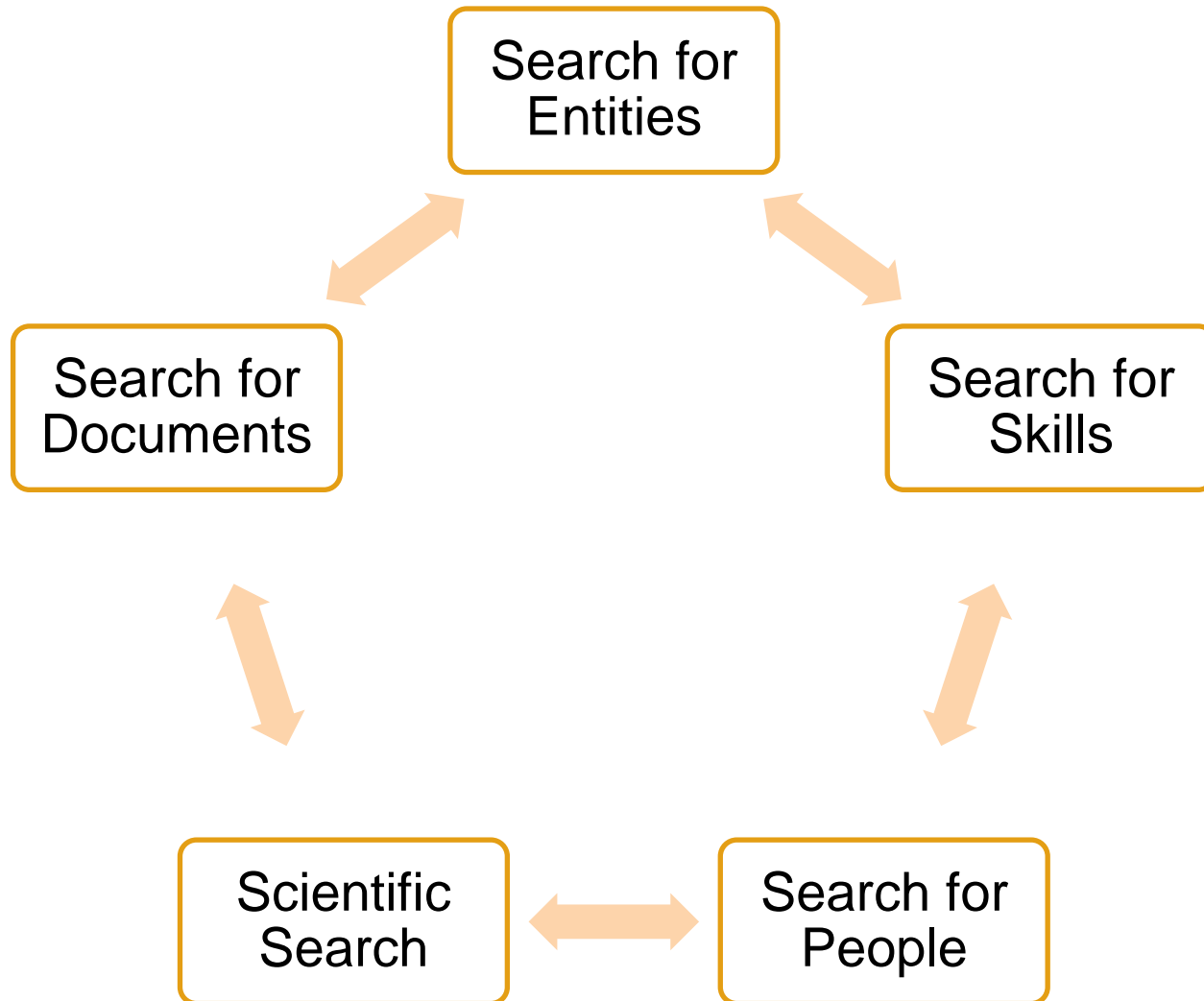
**Submitter**

**Content Owner**

The CTMF supports clarification and resolution of the term meaning, until a Metastore ID can be assigned to the request: e.g. MVNTMPXXXX is a synonym of known concept MVNTAX1234

At any point, the application can check the state of the temporary ID, and update its information with the correct Metastore ID

Once resolved, information will appear in the Metastore. The shortest possible update follows a weekly cycle.

NOVARTIS

*The CTMF in perspective*

- The CTMF has deep implications in the way "concepts" and terminologies are used and produced. As such it is a system that needs to evolve with its user base.

- For this reason, we have released a first core-set of functionalities.

# Search

Search for Entities

Search for Documents

Search for Skills

Scientific Search

Search for People

# Search for entities

- Genes, proteins, target classes, pathways

- Cell lines, cell types, diseases, anatomy

- Compounds, ADC

- Projects...

# Search for Skills

- Disciplines, Lab skills, Instruments, Methods, Technologies, Assays, ...

- Tools/applications

# Search for People

- People name, Department, Sub-department, Location

# Scientific Search

- **Search for**
  - assays and assay results
  - samples, studies, experiments (ELNs)
  - pre-clinical and clinical data, safety data
  - Studies
  - Clinical trials
  - Project information, project decisions
  - Complex queries

# Search for documents ...

- Search for documents in SharePoint/Intranet, Document Databases

- Search for Publications, Patents

- Search for Applications, Tools

- Search for training courses

- ...

# Contribution to scientific queries

NIBR Ontologies can be used for multiple purposes e.g.


Registration systems (metadata capture),  data curation


Data  mapping, bridging ontologies, navigation between concepts and referential data


Text and data mining


Contribution to scientific queries, semantic data federation

U NOVARTIS

# Querying Metastore and other data sources

SPARQL queries

Querying Materialized Views

Federated SPARQL queries with other triplestores

Combining queries in e.g. a datasource, a DWH with querying Metastore materialized views

NOVARTIS

# SPARQL Query

# SPARQL Query

- Give me all PRODUCT combos for component(s) acting as "Renin inhibitor" (MOA) and its subclasses.

# Acknowledgements

- Pierre Parisot

- Andrea Splendiani

- Katia Vella

- Daniel Cronenberger

- TMS, Data Curation team

NOVARTIS