



**C<sup>5</sup>I CENTER**



# System Architecture for Recombinant AI (SARAI)

AI4SE & SE4AI Research and  
Application Workshop 2024



**Dr. Joshua Poore**

Associate Research Scientist with the Applied Research Laboratory for Intelligence and Security (ARLIS)



**Saurabh Srivastava**

PhD Student



**Nikita Patil**

MS Student



**Gaurav Singh**

MS Student



Undergrad. Student Annotators



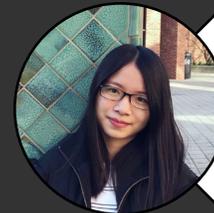
**Dr. Ali Raz**

Assist. Prof. Systems Engineering and Operations Research  
(Expertise: Complex Systems and System of Systems)



**Dr. Paulo Costa**

**Professor – Cyber Security Engineering**  
(Expertise: Information Fusion and Cyber Security)



**Dr. Ziyu Yao**

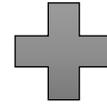
Assist. Prof. Computer Science  
(Expertise: Natural Language Processing)



**Dr. Shou Matsumoto**

Research Assist. Prof. C4I and Cyber Center  
(Expertise: Probabilistic Programming and Ontology)

# System Architecture for Recombinant AI (SARAI) Objectives

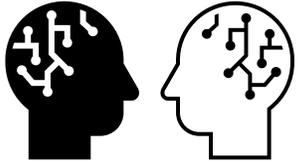


Provide reusable technical assets to expedite and inspire future research into SE4AI

- Quick prototyping, combination, and reuse of data and AI software assets

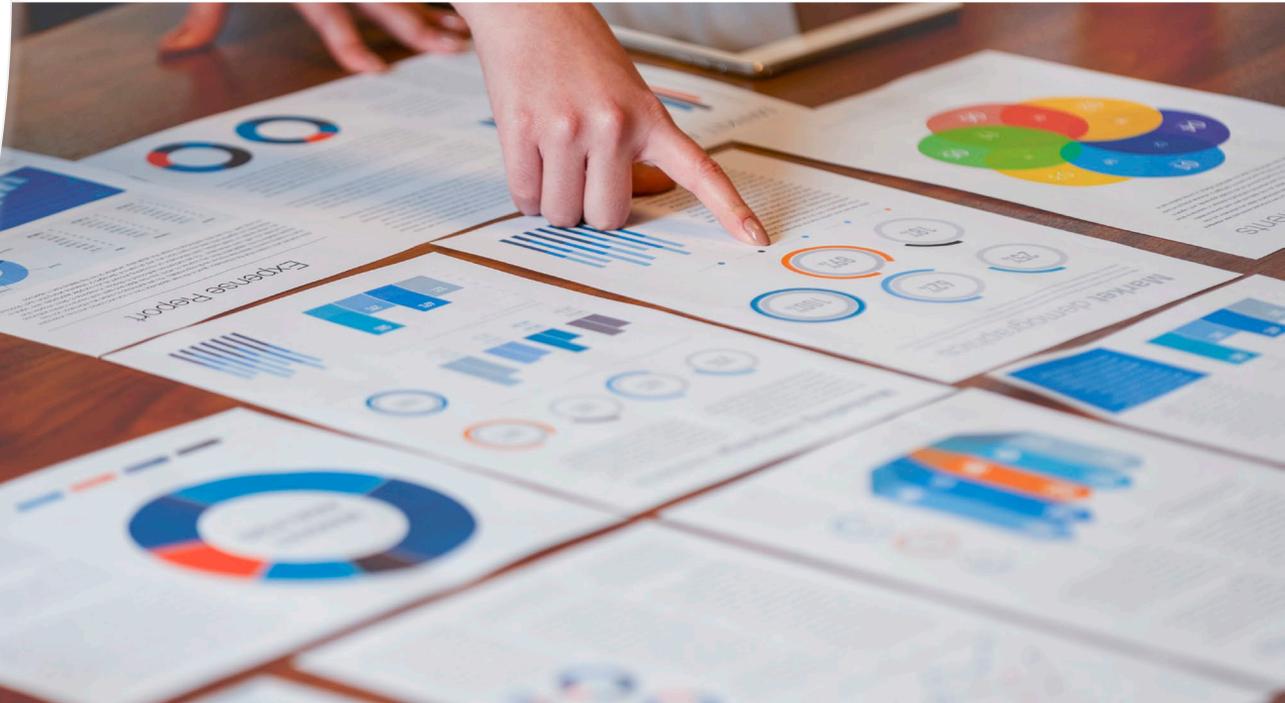
Proof-of-concept experimentation to demonstrate critical capabilities that shape how AI is acquired and evaluated

- Address the reuse, repurpose, obsolescence of current and historical AI applications
- acquisitions, machine learning, autonomy, advanced analytics...



## SARAI Tasks/Works

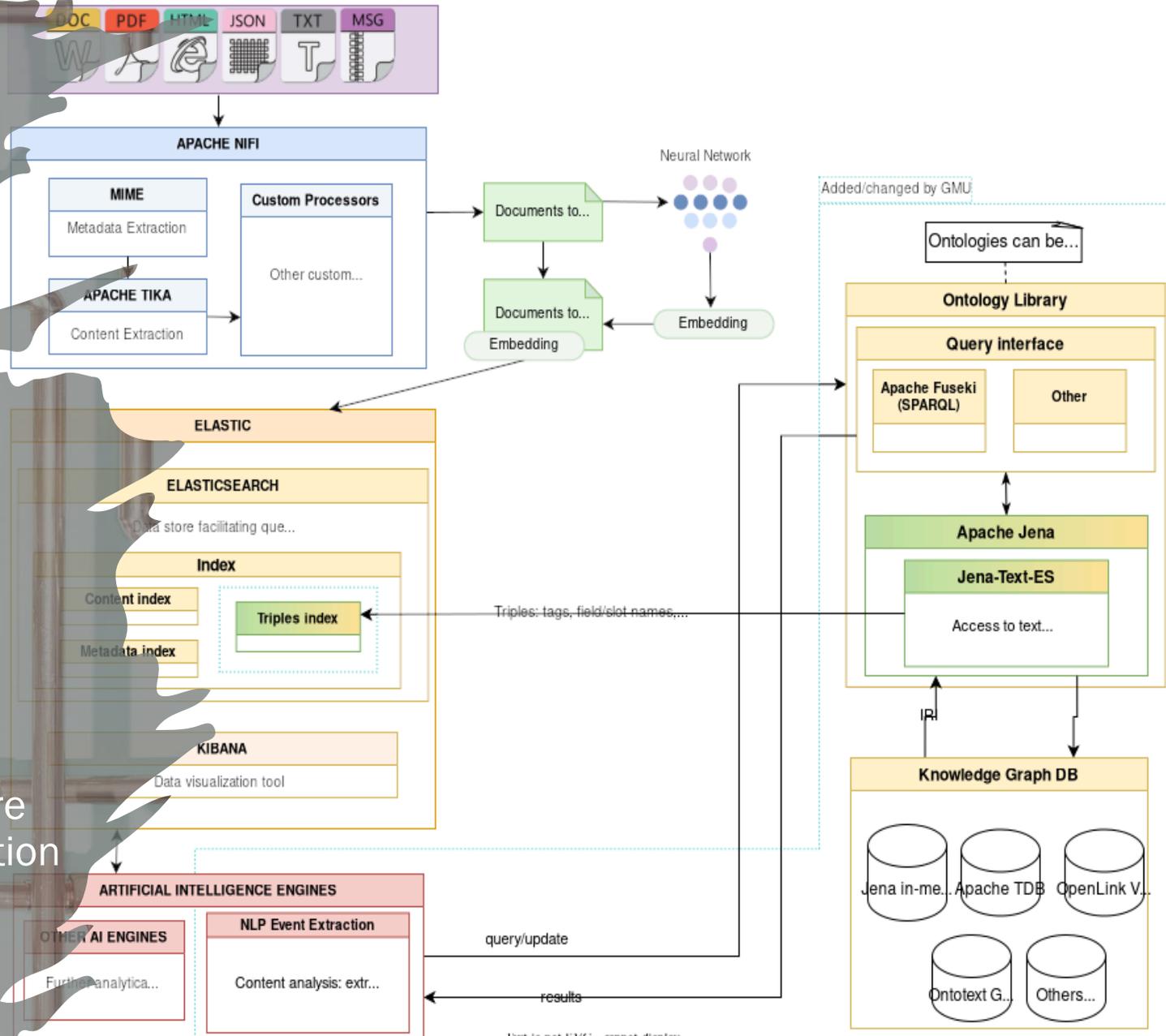
- GMU/C5I center worked towards:
  - development and integration of a pipeline of AI solutions within the scope of document processing and triage.
- This work exercises/provides:
  - basis for demonstrating how parallel lines of AI development efforts and legacy works can benefit from underlying system architecture
    - Methodically integrated various AI solutions



# Methodology

## Synopsis:

Our recombinant AI engineering pipeline develops a system engineering architecture with a focused ontology to provide integration of separately developed AI tools...



# Event Extraction form Ontological Framework and NLP

- **Objective:**

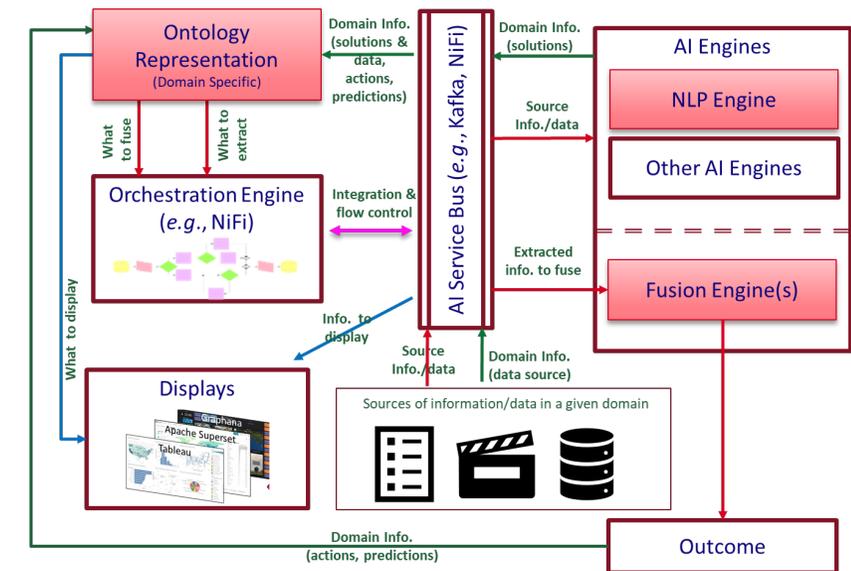
- Complement data ingest pipeline to enable event extraction via an ontological framework

- **Solution:**

- Combine Natural Language Processing (NLP) with a domain-specific Ontology to identify key entities of interest and how these entities are interlinked

## Discriminating Features

- Knowledge base construction and population via Machine Learning
- Semi-automatic ontology development (Machine Learning + “curators”)
- (Natural language-based) user interface
- Fact verification & aggregation
- Domain-specific applications, such as email triage



## Core Dependencies

### Ontology:

- [Apache Jena, Fuseki2, and TDB2](#)
- [Protégé](#) (GUI)
- [OWLGrEd](#) (data visualization)

### NLP tools:

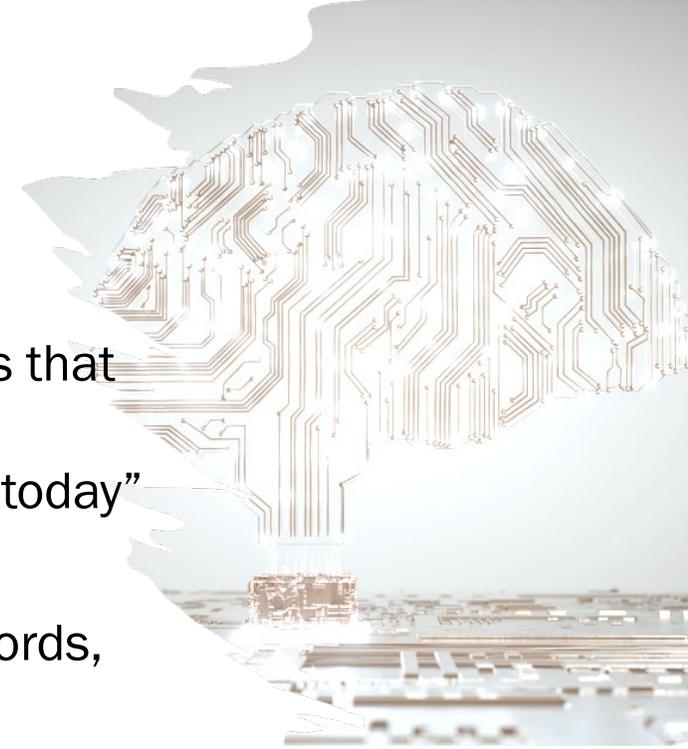
- [Hugging Face Hub/API/Widget](#)
- [PyTorch](#)

### Probabilistic reasoning (future):

- [UnBBayes framework](#)

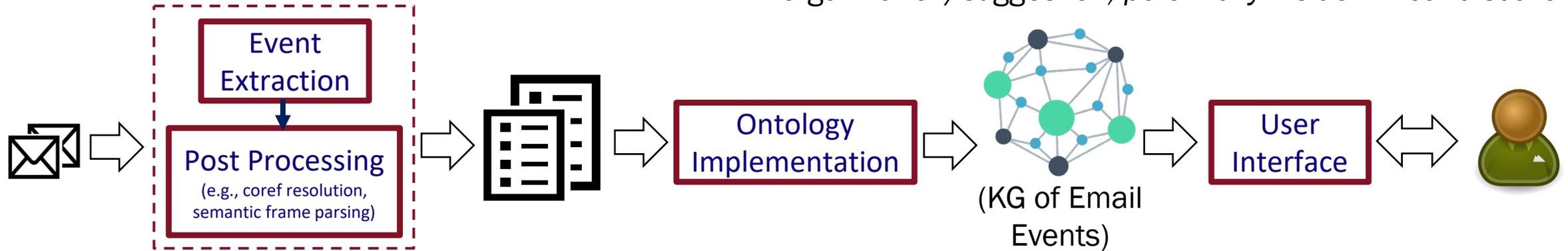
# Neural Network NLP Modules for Email Event Extraction

- A three-stage pipeline
  - 1) Trigger identification
    - Given an email thread, extract words that signal the types of event
    - E.g., “please send me the summary today”
  - 2) Event type prediction
    - Given an email thread and trigger words, predict the type of event
    - E.g., “send me the summary” → Request Data
  - 3) Argument extraction
    - Given an email thread, the trigger words, and the event type, extract word spans as arguments/roles for the event
    - E.g., “the summary” is the requested data, “today” is the requested date



# Instantiation: Constructing a Domain-specific Knowledge Graph

Domain: Email data (for email triage, intelligent task organization/suggestion, potentially insider threat discovery, etc.)



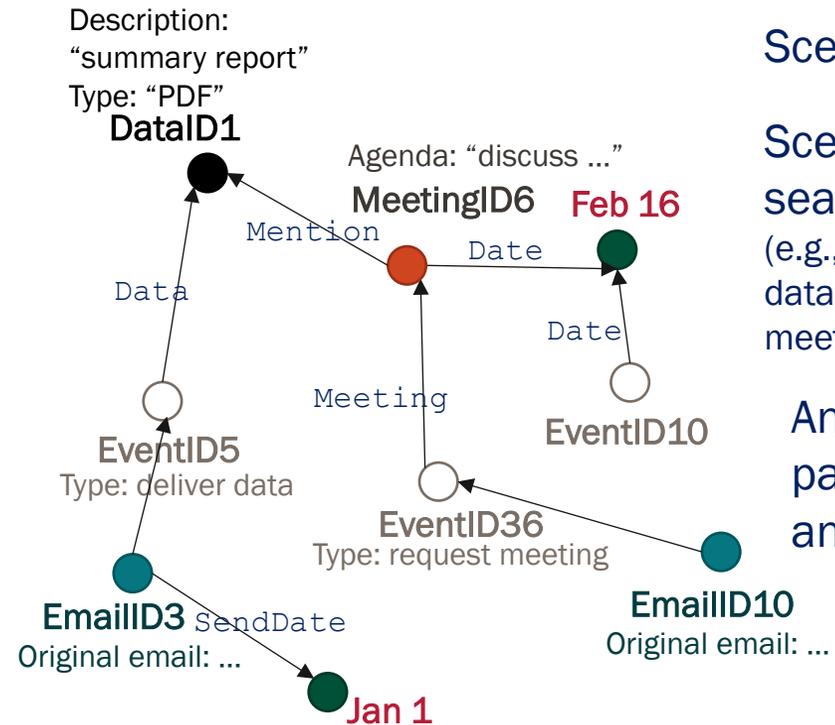
**EmailID3**  
 Hi Andrew, the attached is a PDF of the summary report -- Alice

**EventID5**  
 EventType: Deliver Data  
 Data: **DataID1**  
 |- Type: "PDF"  
 |- Description: "summary report"

**EmailID10**  
 How about meeting this Wed to discuss the report from Alice?

**EventID36**  
 EventType: Request Meeting  
 Meeting: **MeetingID6**  
 |- Agenda: "discuss the report from Alice" → **DataID1**  
 |- Date: "this Wed"

Feb 16



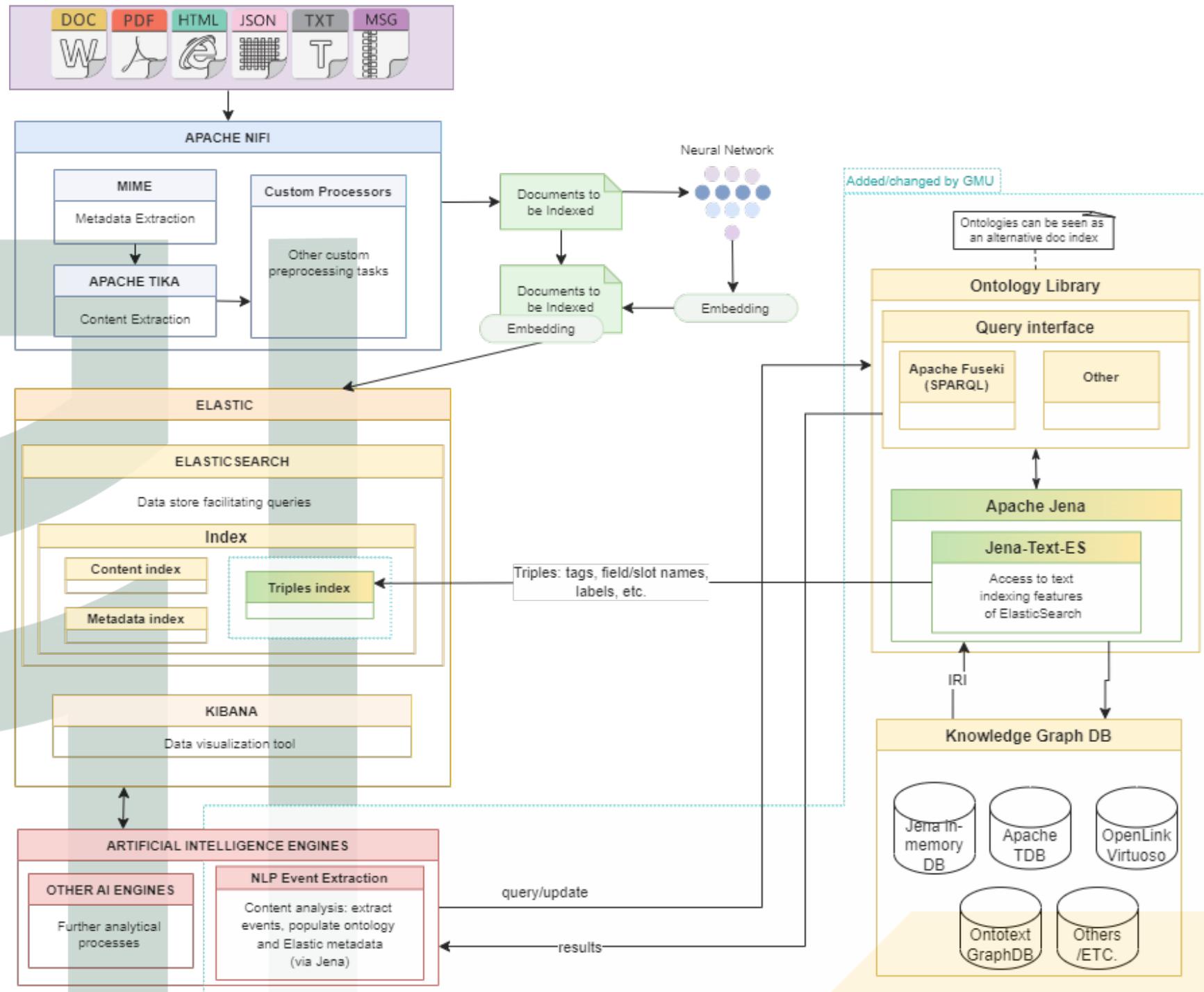
Scenario 1: SPARQL query

Scenario 2: Semantic search via natural language (e.g., "show me a description of the data we will discuss in the Feb 16 meeting")

And more: reasoning, pattern learning for anomaly detection, etc.

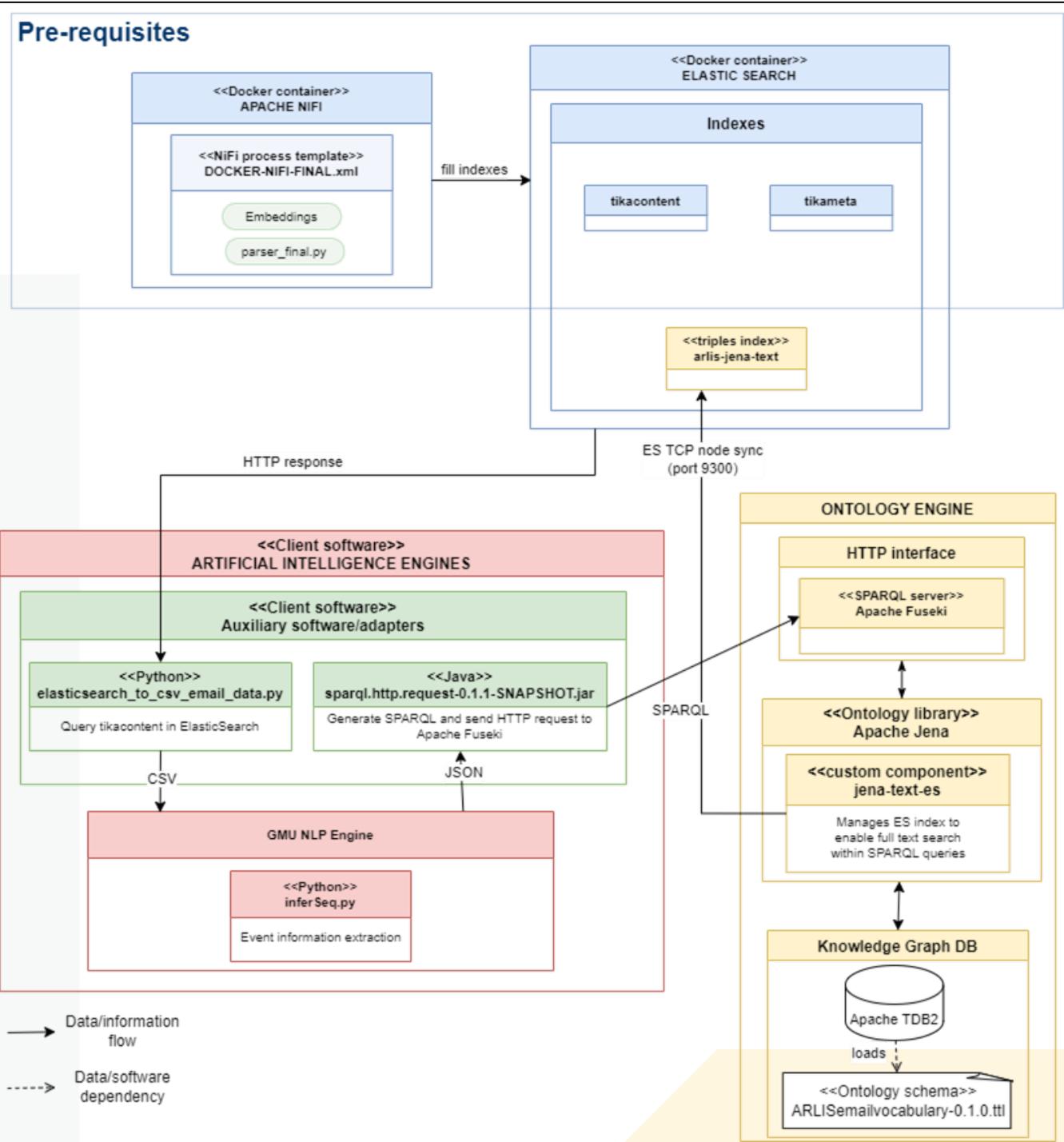
# GMU-ARLIS System Architecture

Components' collaboration model



# Prototype/Demo Implementation

The top screenshot shows the Apache Fuseki web interface. It displays a SPARQL query editor with a query for finding meeting events. The query results are shown in a table with columns for email, event, meetingName, location, members, meetingDate, and meetingagenda. The bottom screenshot shows the OWLGrEd ontology visualization tool, displaying a complex network of classes and instances related to meetings and events.





# Prototype/demo implementation – classification/inference

Simplified/filtered for better visualization. Names are redacted.

## NLP event classification/extraction

```
{
  "ID":
  "b56021de4258c3cebf37346b4bc9876e1944a9c7c7d7
  aebf3ec79487f402b6c6",
  "Request_Meeting": [ {
    "Meeting Name": "meeting",
    "Meeting Date": "2022-03-23T11:15:00",
    "Meeting Location": "room 3127",
    "Meeting Members": [ "████", "████" ],
    "Meeting Agenda": "introduce █████ and
    get some input on various projects she will
    be assisting us with"
  } ]
}
```

Some common concepts were reused from BFO/CCO

- Basic Formal Ontology (BFO): <https://basic-formal-ontology.org/>
- Common Core Ontologies (CCO): [github.com/CommonCoreOntology](https://github.com/CommonCoreOntology)

## Sample SPARQL: find members with conflicting times

Dataset: /ARLIS

query upload files edit info

SPARQL query

To try out some SPARQL queries against the selected dataset, enter your query here.

EXAMPLE QUERIES

Selection of triples Selection of classes

PREFIXES

rdf rdfs owl xsd

SPARQL ENDPOINT /ARLIS/query

CONTENT TYPE (SELECT) JSON

CONTENT TYPE (GRAPH) Turtle

```
5 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
6 PREFIX arlis: <https://c4i.gmu.edu/arlis/ontologies/ARLISemailvocabulary.ttl#>
7 PREFIX : <.>
8
9 SELECT distinct ?individual WHERE {
10   ?individual a arlis:Person .
11   ?event1 a arlis:Event .
12   ?event1 arlis:hasMeetingMembers ?individual.
13   ?event2 a arlis:Event .
14   ?event2 arlis:hasMeetingMembers ?individual.
15   FILTER (?event1 != ?event2)
16 }
17 ORDER BY ?individual
```

QUERY RESULTS

Table Raw Response

Showing 1 to 1 of 1 entries

Search: Show 50 entries

individual
1 <http://www.c5i.gmu.edu/arlis/ARLISemailDataSamples/person1>

Showing 1 to 1 of 1 entries

# Prototype/demo implementation – visualization tools

The screenshot shows the OWLGrEd web application interface. At the top, the browser address bar shows the URL `owlgred.lumii.lv/online_visualization/4e1a`. The navigation menu includes **OWLGrEd**, Home, Get Started, Visualize Online (active), Notation, Extensions, R&D, Success Stories, Users, and About Us.

The main content area features the heading "Try our ontology visualization" and a sub-heading "Look up our graphical notation for ontologies." Below this, it states: "Currently best results with moderate size ontologies. Diagram will show only the direct ontology." A call to action asks users to share impressions and suggestions to [owlgred@lumii.lv](mailto:owlgred@lumii.lv).

On the right side, there are three buttons: "Enjoy our examples" (with an eye icon), "Visualize your ontology" (with a document icon), and "Share the link" (with a share icon). Below these buttons is a text input field containing the URL `http://www.semanticweb.org/owl/owlapi/turtle#` and a "turtle.png" button.

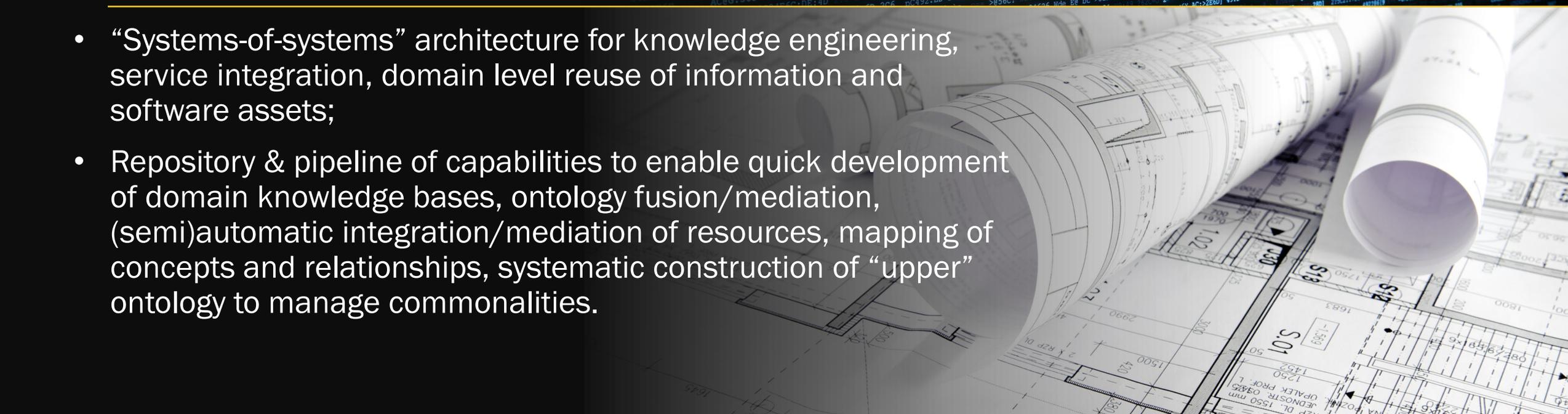
The central part of the interface displays a complex ontology visualization. It consists of a large grid with various colored nodes (yellow, green, grey) and connecting lines. A vertical toolbar on the left side of the visualization area includes a zoom-in (+) button, a zoom-out (-) button, and a refresh/clear button.



# Results/outcomes

## Synopsis:

- “Systems-of-systems” architecture for knowledge engineering, service integration, domain level reuse of information and software assets;
- Repository & pipeline of capabilities to enable quick development of domain knowledge bases, ontology fusion/mediation, (semi)automatic integration/mediation of resources, mapping of concepts and relationships, systematic construction of “upper” ontology to manage commonalities.



# Conclusion and Key Contributions

- **System Architecture & ontology** development
  - Use cases of systems integration & semantic technologies
  - Assists reuse, prevents obsolescence of AI solutions
  - Agile/faster prototyping, development, V&V of AI pipelines
- **NLP components** that produce ontological **knowledge base** for **document triage**
  - Framework for parallel data processing, enterprise level data indexing, data curation...
- Technical products disseminated as **open-source** assets
  - Assist in scaffolding technically similar efforts
  - Provide AI Engineering examples for a larger community



# Future Research



## Domain adaptation

- Currently, information extraction is on email data only
- Many other domains: science, medicine, cybersecurity, etc.
- Can we build a similar system using as few annotations as possible?



## Data fusion across different knowledge sources

- Emails, X (social media), news articles, policy documents...
- Text, tables, images, etc.



## Language interface to intelligence assistants

- Interpretability & interaction
- Analyst Query



C<sup>5</sup>I CENTER

Thank you for your  
attention

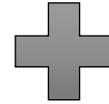




C<sup>5</sup>I CENTER

Backup slides

# Recombinant AI Seedling Objectives



- Perform proof of concept experimentation to demonstrate critical capabilities that will shape how artificial intelligence is acquired and evaluated;
- Provide reusable technical assets to expedite and inspire future research into artificial intelligence;
  - Explore reusable and recombinant AI/ML through flexible data engineering pipelines and efficient Data Service architectures
- Provide examples of engineering methodology and process to scaffold the AI Engineering discipline and develop the ODNI's current and future workforce.

# GMU Team Sow/Tasks

- Develop analytical frameworks (e.g., ontological framework) to enable heterogeneous data exploitation and fusion in support of Recombinant AI objectives.
- Building on the outcomes of Natural Language Processing (NLP) applied to various documents in a given domain (e.g., data extraction, indexing, and translation etc.), an ontology framework, for example, will identify key entities of interest in that domain and how these entities are interlinked towards inferring root causes or potential future courses of action.
- Such frameworks set foundations for exploiting and fusing heterogeneous data to overcome limitations of missing data and/or extract new information from disparate and siloed data sets.



# Dataset: MailEx

- A conversational email-domain event extraction dataset
  - Source: the Enron email corpus
  - ~1,200 email threads
  - ~3,400 individual emails
  - ~4,500 events
  - (numbers after handling annotator disagreement)
- 11 event types (with arguments)

Event Type	%	Frequent Argument Roles
Request Data	8.91	Data IdString (72%), Request Members (23%), Request Date (2%)
Request Action	20.22	Action Description (54%), Action Members (35%), Action Date (6%)
Request Meeting	5.02	Meeting Members (31%), Meeting Agenda (21%), Meeting Date (18%)
Request Action Data	2.39	Action Description (51%), Action Members (38%), Request Members (8%)
Request Meeting Data	0.71	Meeting Members (31%), Meeting Agenda (21%), Meeting Date (18%)
Deliver Data	24.32	Data IdString (48%), Data Value (39%), Deliver Members (10%)
Deliver Action Data	28.72	Action Description (46%), Action Members (41%), Action Date (9%)
Deliver Meeting Data	6.21	Meeting Members (34%), Meeting Date (19%), Meeting Time (12%)
Amend Data	2.22	Amend Members (26%), (Context) Data IdString (25%), (Revision) Data Value (25%)
Amend Meeting Data	1.27	(Revision) Meeting Time (22%), (Revision) Meeting Date (19%), (Context) Meeting Name (16%)

Table 3: Distributions of event types (in percentage) and frequent argument roles in MAILEX.

# Intelligent ML/NLP Component

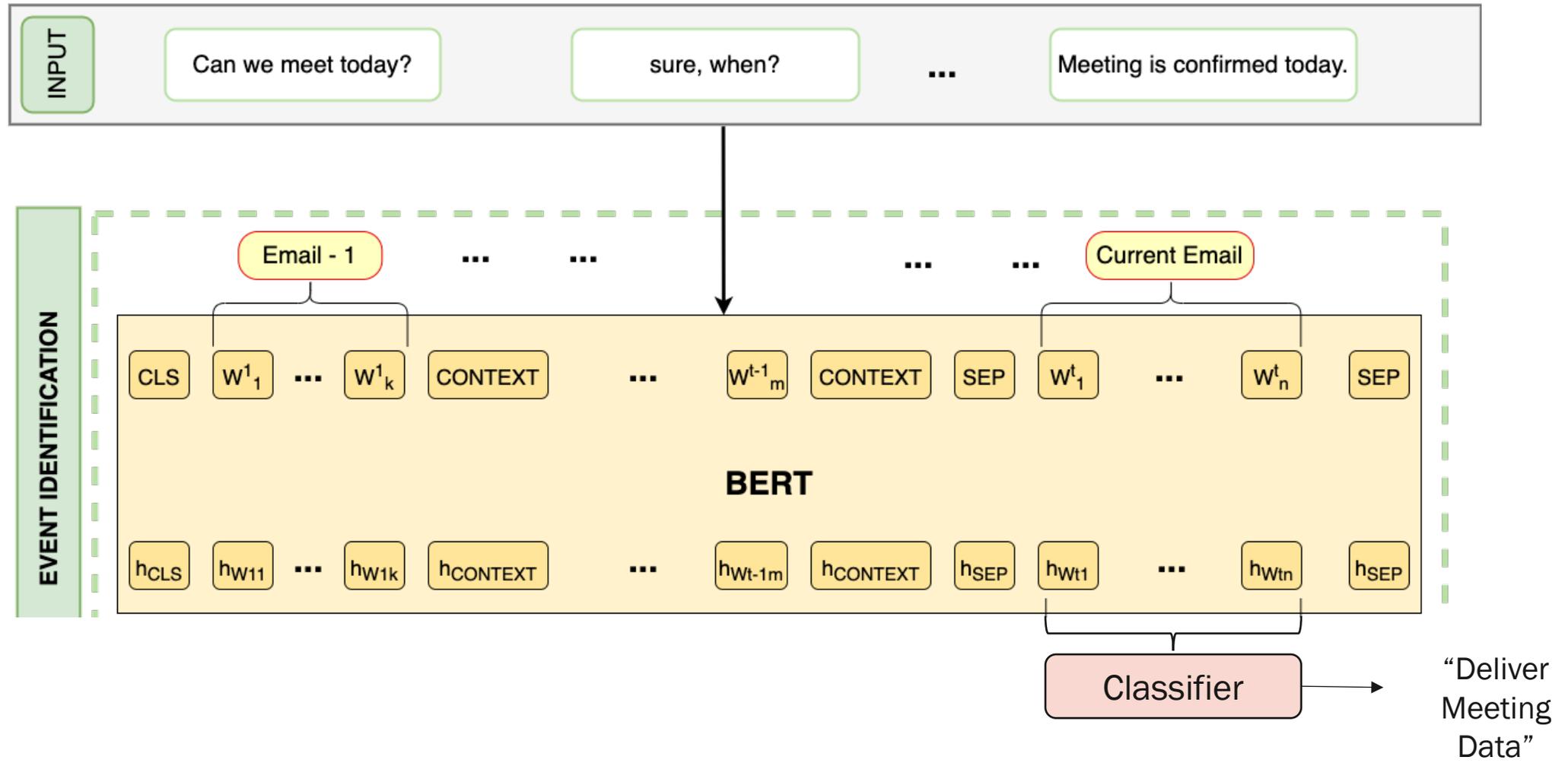
- A neural network model for event extraction on email threads



The screenshot shows the 'Enron Annotation Interface' with navigation links for 'Annotation Guidelines' and 'Interface Manual'. It features a 'Light Theme' toggle and a 'Logout' button. A 'Total Events Added' counter is shown, along with a timer at '00:00:29'. The interface is currently on 'Current Turn: 2', with 'From' and 'To' dropdown menus and a 'Show Turn' button. A 'Not Sure?' button is also present. Below this, there are buttons for 'Show Turn: 1', 'Select Event-Type', 'Finish Current Event', and 'Submit'. The main content area displays a grid of words from an email thread, each with a dropdown menu for annotation. The words shown are: Paul, -, Did, you, ask, the, Market, Services, rep, to, credit, the, commodity, on, PNM, 's, bill, and ?. The dropdown menus are currently set to '0'.

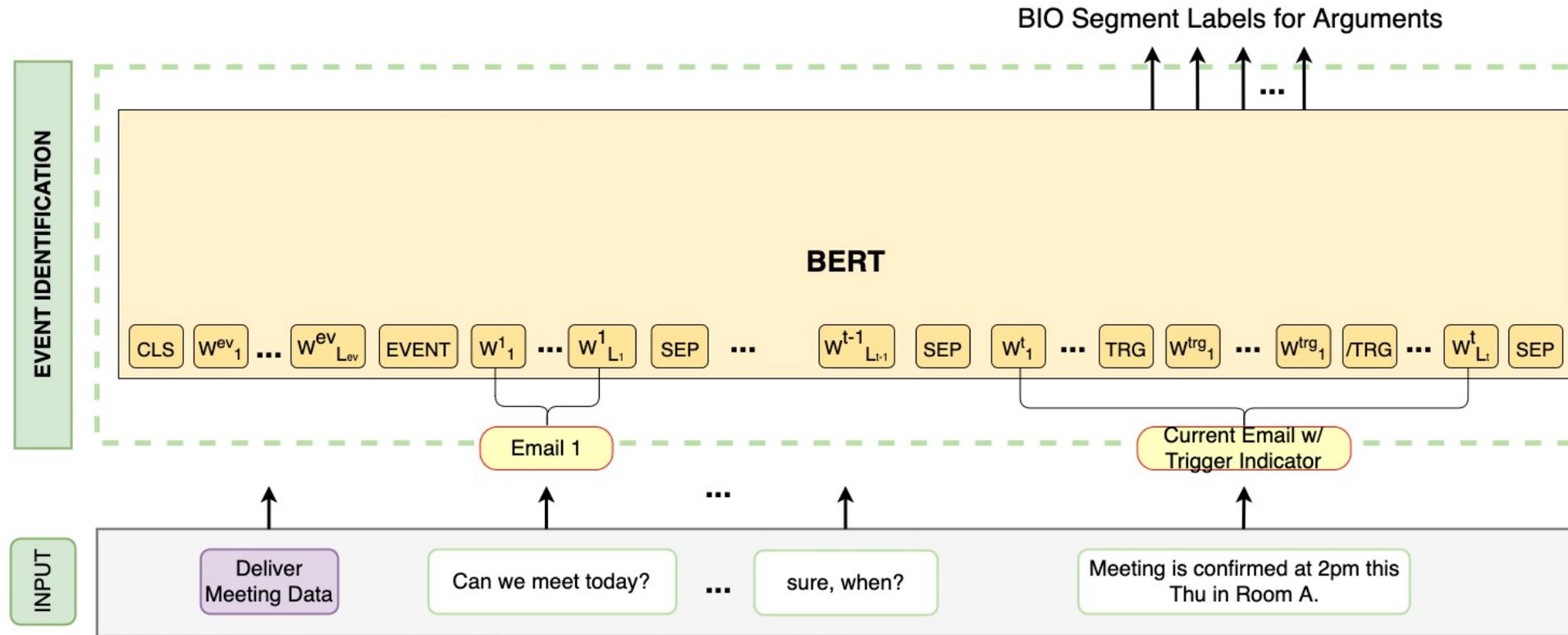
# Neural Network Modules for Email Event Extraction

- Event Identification
  - Given the current email and the thread history, identify event types in the email



# Neural Network Modules for Email Event Extraction

- Argument Extraction

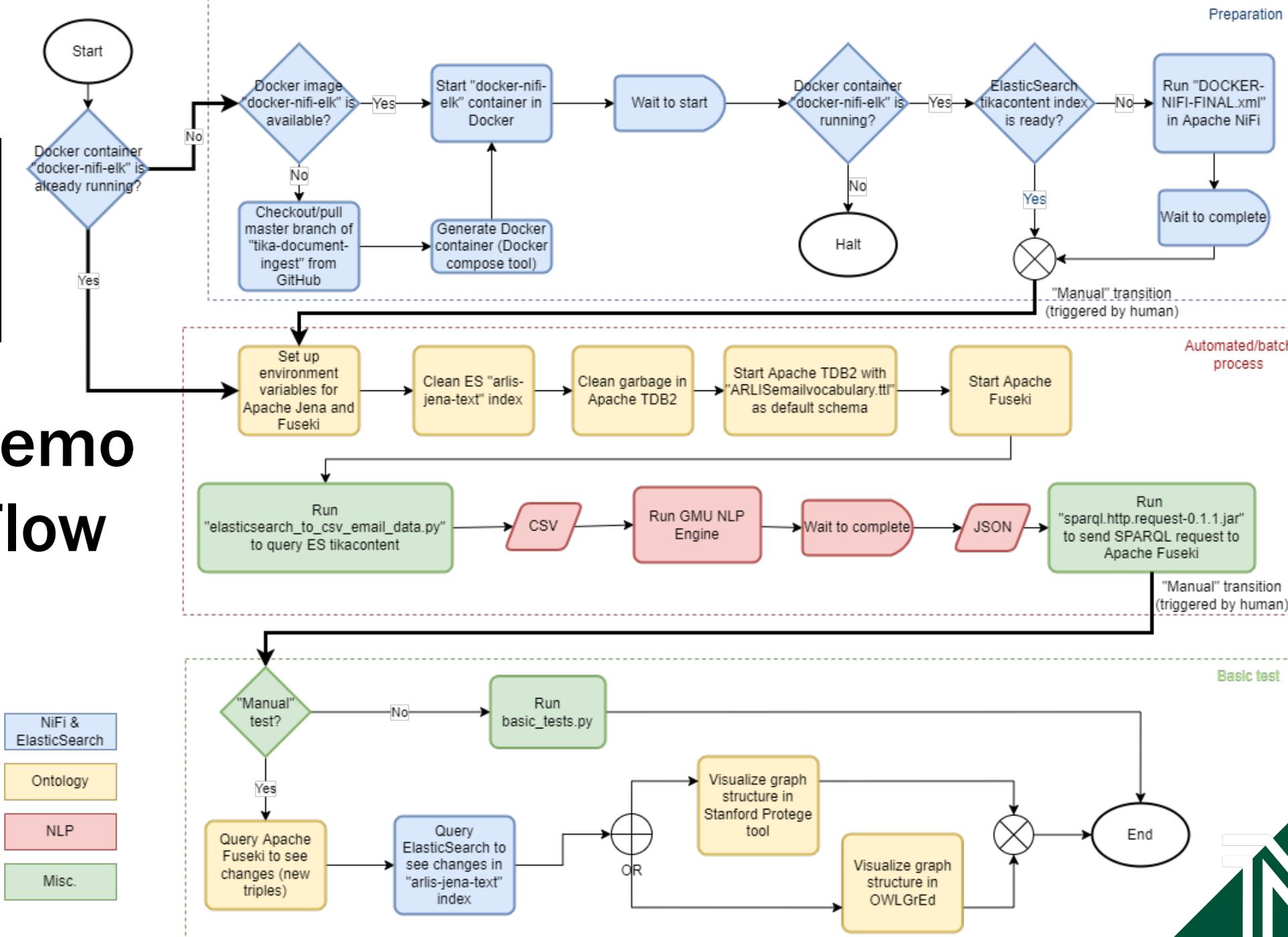
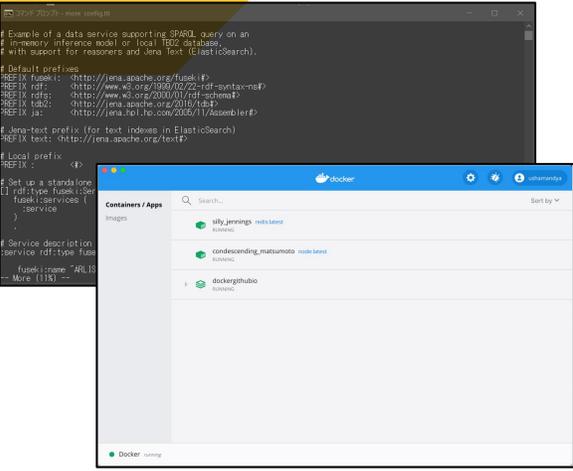


e.g.,  
[B-DeliverMeet:Date,  
I-DeliverMeet:Date]  
gives a "Date" span



Deliver Meeting:  
- Time: 2pm  
- Date: this Thu  
- Location: Room A

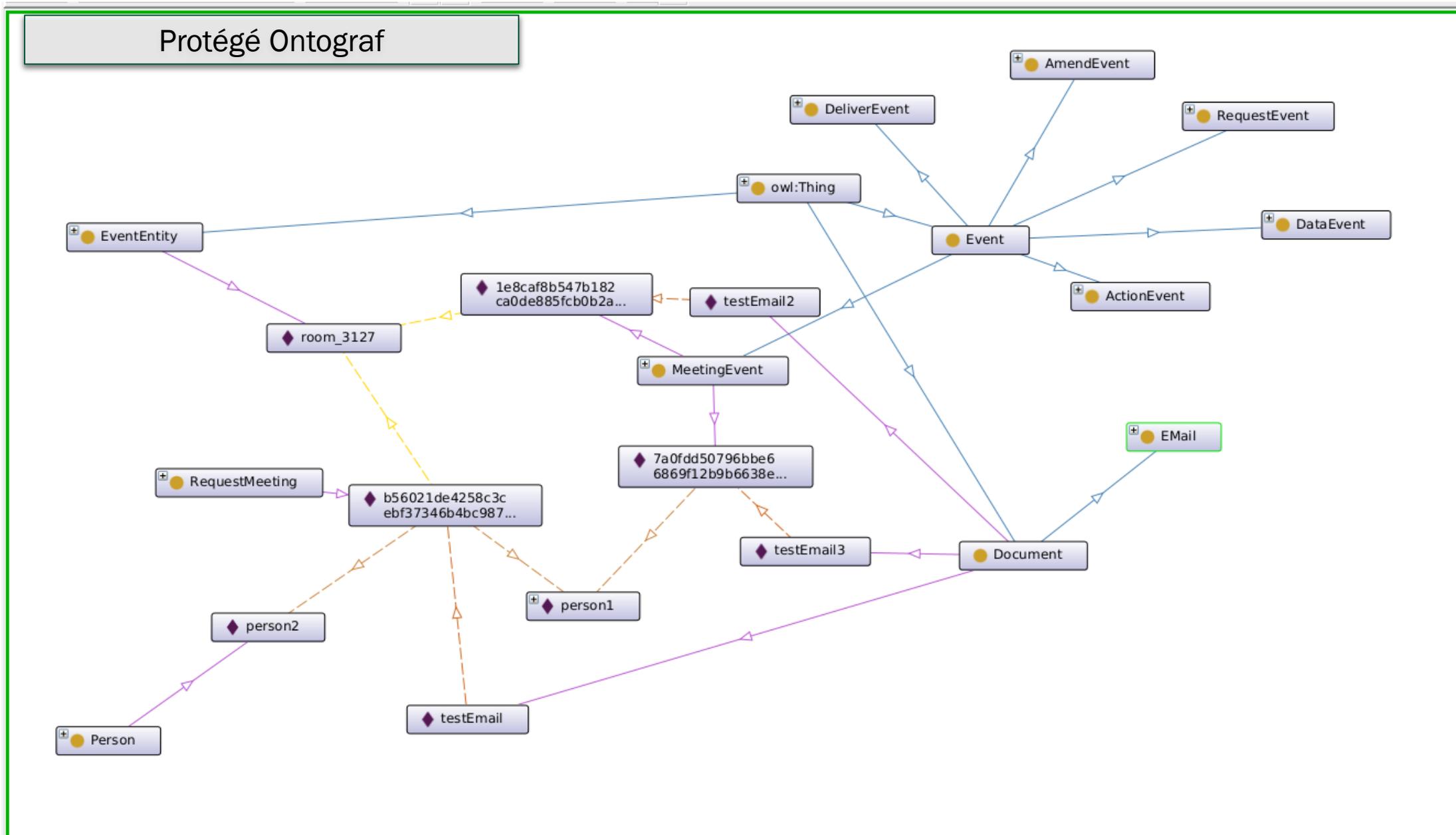
# Prototype/Demo Execution Flow



- NiFi & ElasticSearch
- Ontology
- NLP
- Misc.



# Prototype/demo implementation – visualization tools



# Domain Adaptation of Event Extraction

- Research Objective: Enabling fast development of domain-specific KG from data
  - E.g., adapting from open domain to emails (our own dataset), AI Incident Database [1], COVID tweets [2], cybersecurity documents [3,4], BioNLP [5], scientific research [6]
  - Some topics to explore:
    - Pre-training, data augmentation, etc.
    - Uncertainty
      - Analyze and explore/exploit epistemic (ignorance, lack of data, low comprehensiveness) VS aleatoric (ubiquitous, natural in domain) uncertainty.
        - Outlier, out-of-distribution analysis, etc.
      - Identify events/information summaries to report as AI Incident

[1] <https://incidentdatabase.ai/>

[2] Shi Zong, Ashutosh Baheti, Wei Xu, and Alan Ritter. Extracting a knowledge base of covid-19 events from social media. arXiv preprint arXiv:2006.02567, 2020.

[2] Taneeya Satyapanich, Francis Ferraro, and Tim Finin. Casie: Extracting cybersecurity event information from text. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pages 8749– 8757, 2020.

[3] Hieu Man Duc Trong, Duc-Trong Le, Amir Pouran Ben Veyseh, Thut Nguyn, and Thien Huu Nguyen. Introducing a new dataset for event detection in cybersecurity texts. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5381–5390, 2020.

[4] Jin-Dong Kim, Tomoko Ohta, Sampo Pyysalo, Yoshinobu Kano, and Jun'ichi Tsujii. Overview of bionlp'09 shared task on event extraction. In Proceedings of the BioNLP 2009 workshop companion volume for shared task, pages 1–9, 2009.

[5] Luan, Yi, Mari Ostendorf, and Hannaneh Hajishirzi. "Scientific Information Extraction with Semi-supervised Neural Tagging." Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing. 2017.

# Cross-domain knowledge/data fusion

- Research Objective: Fusing the cross-domain knowledge/data for applications
  - Fusion across modalities for higher-level inference
    - Knowledge graph + Text (and more in the long term)
    - Applications: recommendation, question answering, search, etc.
  - Potential Application Domains with domain knowledge graphs
    - COVID/new virus discoveries
      - News articles reporting the unusual/novel pathogen activity (e.g., COVID cases)
      - CDC websites describing the disease symptoms and prevention
    - Emerging AI/ML Research
      - AI/ML papers and blogs
      - Social media discussion on AI/ML applications
    - Fact verification/misinformation detection
      - Web-pages, tables, tweets with miss information
  - Techniques:
    - Leveraging our domain adaptation techniques
    - With humans in the loop for reliability
    - Explainable AI techniques, probabilistic ontology/reasoning...
    - ...

