

Need for text mining

Number of new scientific publications is growing rapidly

New terms (genes, proteins, chemical compounds, drugs) are constantly created

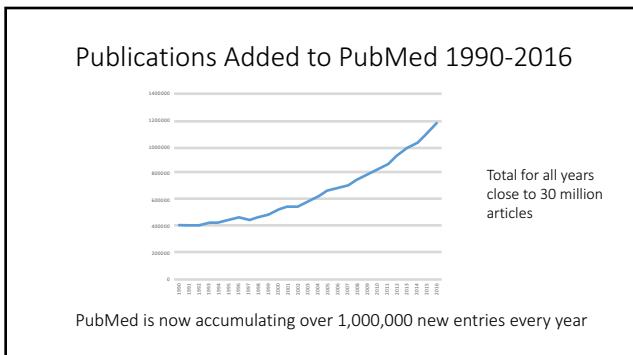
Impossible to manage such an information overload!

Scientific Information Overload

The global research community generates ~2.5 million new scholarly papers per year (English only)
A new research paper is published every 12 seconds
70,000 papers published on a single protein

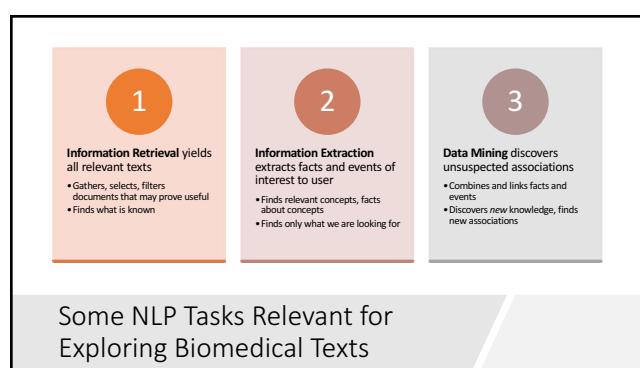
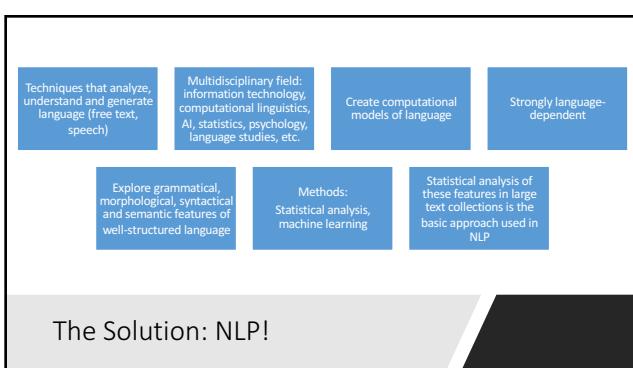
Challenge to scientists

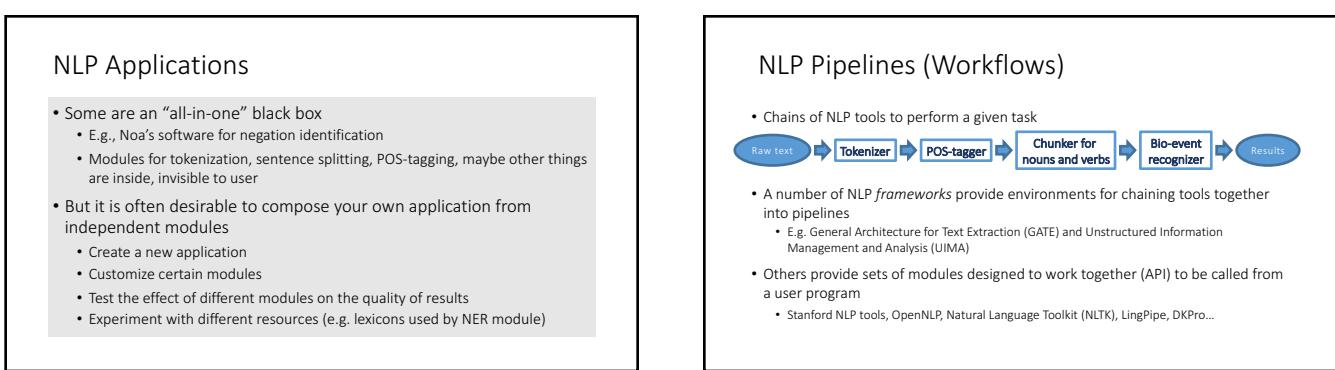
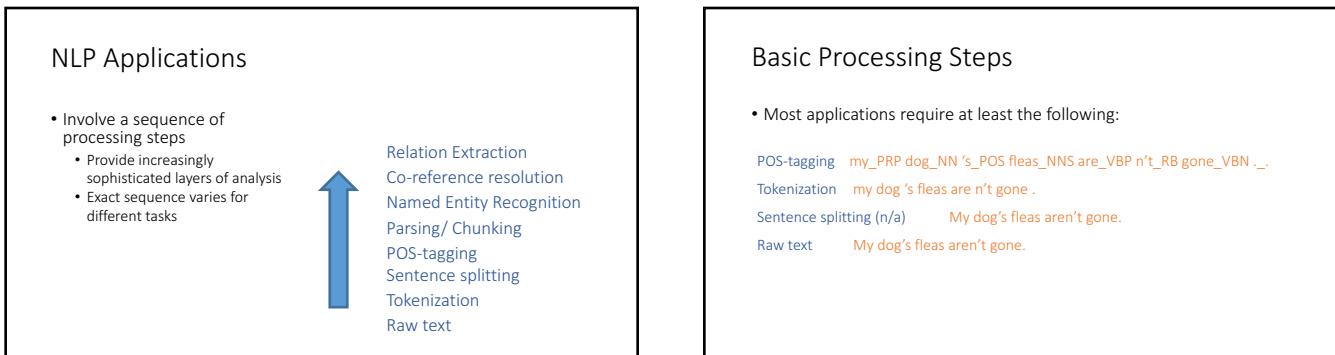
Keep updated on new developments, paper writing, project proposal preparation, paper reviewing, peer assessment, etc.



| Need for text mining

- Increased availability of full text – Information retrieval is an insufficient solution
- Bio-databases, controlled vocabularies and bio-ontologies encode only small fraction of information
- Most information is in textual form – unstructured data
- Automated aids are needed!





Typical Scenario

- A scientist wants to apply text mining techniques to find articles including references to certain entities (e.g., proteins, genes) and their interactions
 - Knows nothing about NLP or Computer Science
 - Unfamiliar with NLP technologies
- Searches for NLP software that might help

Typical Scenario

- Finds existing tools and frameworks that are freely available

Not to mention several commercial (i.e., pricey) options


- Questions
 - Do these things all do the same thing, or do they differ in some way?
 - Do some work better than others?
 - Are some easier to use than others?
 - How does one choose?



Issues

- Most of these tools provide general NLP support, not geared to BioNLP
- Why does this matter?

Characteristics of biomedical literature

1

Heavy use of domain specific terminology (12% biochemistry related technical terms)

- Examples: chemoattractant, fibroblasts, angiogenesis

2

Constant introduction of new terms and short forms or abbreviations

Characteristics of biomedical literature

3

Polysemic words

Examples: APC: (1) Argon Plasma Coagulation (2) Activated Protein C; or teashirt: (1) a type of cloth (2) a gene name (tsh).

4

Heavy use of acronyms

Examples: Activated protein C (APC), or vascular endothelial growth factor (VEGF)

Characteristics of biomedical literature

5a

Term Ambiguity

Gene terms may be also common English words

- BAD human gene encoding BCL-2 family of proteins (*bad news, bad prediction*)

5b

Term Ambiguity

Gene names are often used to denote gene products (proteins)

- *suppressor of sable* is used ambiguously to refer to either genes and proteins

Characteristics of biomedical literature

6

Typographical variants (e.g. in writing gene names)

- Example: TNF-alpha and TNF alpha (without hyphen)

7

Most words have low frequency (data sparseness)

No Problem!

Lists of NLP Tools for Biomedical texts

- <http://bionlp.org>
- http://biocreative.sourceforge.net/bionlp_tools_links.html
- <http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=7784950>
- <http://www.nactem.ac.uk/index.php> (but also lots of other things)

However...

- Many of these tools are difficult to install, configure, and use without some computational expertise
- Even more difficult to modify or adapt without computational expertise and some knowledge of NLP
- Also: which tools performing the same task perform best and/or are best suited to a given task?

Another Sneaky Underlying Problem

- Input and output of tools from different sources differ dramatically!!!
- Often demands significant effort and expertise to adapt tools from different sources to work together
...if it is possible at all
- I.e., tools are not **interoperable**

Lack of Interoperability

- ▶ Software waste cycle
 - ▶ Because of the difficulties of adapting software developed at other sites to one's data and/or purpose, researchers often write their own versions
 - ▶ Huge waste of time and effort to reinvent the wheel
- ▶ Increasing availability of open tools and data
 - ▶ Open tools and data developed at a single site (e.g. Stanford tools) usable together with no modification
 - ▶ BUT
 - ▶ Want to be able to use tools from different sites together with minimal or no effort to adapt
 - ▶ Important for Open Advancement approach to development
 - ▶ Want to be able to process data in various formats with minimal or no effort to adapt

What is Needed

1

Develop/ provide access to a range of freely available advanced text mining tools specially tailored to scientific publications

2

Enable scientists to easily use these tools without having to be a computer scientist or an expert in NLP

3

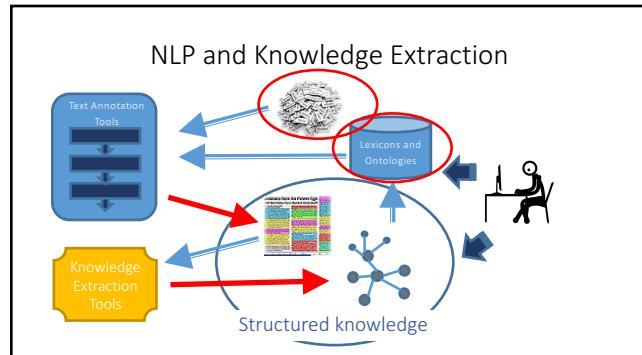
Enable scientist to easily adapt existing solutions to specific domains or problems without having to be a computer scientist or an expert in NLP

Interoperability is key!

Tools are Not Enough

BioNLP also needs **language resources**:

- Large bodies of scientific publications that can be searched and mined for information and knowledge
 - Large bodies of **annotated** scientific publications that can be used to develop language models (e.g. via machine learning)
 - Lexicons and Ontologies of biomedical terms to assist in entity recognition etc.



Biomedical Literature Repositories

- Electronically accessible to the public through the web
 - Centralized institutional (for example, PubMed) or academic (for instance, Highwire Press and Hollis) repositories of peer-reviewed articles or article abstracts.
 - Article collection repositories hosted by publishers (for example, BioMed Central and EMBASE).
 - Online access to indexed scholar articles retrieved through web spiders and crawlers (for example, Google Scholar and Scirus).
 - Basic and Milinovic (2012) provide one overview of world literature databases for the Biomedical field

Despite all this, a corpus resource bottleneck

- 
 - 1 Lack of freely available, large-scale, richly annotated corpora to support
 - Training of ML algorithms
 - Development of computational grammars
 - Evaluation of text mining components
 - 2 Existing corpora in many different formats
 - Often necessary to extract text from PDF (not always accurate!)
 - Formats such as various home-grown XML, JSON, "BIO", ...
 - 3 Annotations may be named and/or structured very differently
 - Annotations for the same phenomenon have to be transduced to be used together
 - Sometimes there is no direct mapping
 - Not interoperable!

Knowledge Resources for BioNLP

1

Lexical and terminological resources

- Lists of terms / lexical entries

2

Ontological resources

- Concepts with definitions, sometimes including synonyms
- Relations between the concepts, typically is-a and has-part

Table 1. Some knowledge sources for biomedical natural language processing.

Informatics for Integrating Biology and the Bedside (B2b2 - https://www.b2b2.org/)	National Center for Biomedical Computing with focus on translational research that facilitates and proves data sets for clinical natural language processing research
Gene Ontology (https://www.geneontology.org/)	Controlled vocabulary with relationships including paronymy and inheritance, designed for describing gene functions, broadly construed
Entrez Gene (https://www.ncbi.nlm.nih.gov/gene/)	Source of gene names, symbols, and synonyms also the source for GeneIDs and SUMMARY fields
PubMed/MEDLINE (https://www.ncbi.nlm.nih.gov/pubmed/)	The National Library of Medicine's database of abstracts of biomedical publications (MEDLINE) and search interface for accessing them (PubMed)
Unified Medical Language System (https://www.nlm.nih.gov/research/umls/)	Large lexical and conceptual resource, including the UMLS Metathesaurus, which aggregates a large number of biomedical and some genomic vocabularies
SWISS-PROT (https://www.uniprot.org/)	Database of information about proteins with literature references, useful as a gold standard
PharmGKB (https://www.pharmgkb.org/)	Database of relationships between a number of clinical, genomic, and other entities with their references. itself is a gold standard
Comparative Toxicogenomics Database (https://ctdbase.org/)	Database of relationships between genes, diseases, and chemicals, with literature references, useful as a gold standard
Various terminological resources, data sources, and gold-standard databases for biomedical natural language processing.	

Cohen KB, Hunter LE (2013) Chapter 16: Text Mining for Translational Bioinformatics. PLOS Computational Biology 9(4): e1003044.
<https://doi.org/10.1371/journal.pcbi.1003044>

Are the existing knowledge resources sufficient for text mining?

NO!

- Limited lexical and terminological coverage of biological sub-domains
- Terminological variation and complexity of names
- Variation occurs in controlled vocabularies and texts but discrepancy between the two
- Exact match methods fail to associate term occurrences in texts with databases

Furthermore...

- The same interoperability problem exists for resources!
 - Different physical formats
 - PDF, XML, plain text...
 - Extraction of text from PDF can be unreliable
 - Different representations for annotations
 - Different physical formats
 - XML, JSON, brackets, BIO

BIONLP2016 Shared Task Data Sets

```
(nmod_ _structures_2 Chromosomal_1)
(nmod_ _testosterone_3 Pseudomonas_4)
(dobj_4_3 testosterone_5)
(nmod_ _structures_2 of_3)
<>|LRB|LRB Chromosomal|JJ|N/N structures|NNS|N
of|IN|(NP|NP|NP|Pseudomonas|NN|N/N
testosterone|NN|N|.|.
<>|NN|NN|N|.|.
```

CDI	199 202	ORGANIZATION
CDI	375 378	ORGANIZATION
CDI	426 429	ORGANIZATION
CDI	501 502	ORGANIZATION
CDI	713 716	ORGANIZATION
CDI	854 867	ORGANIZATION
Olmsted County	1001 1015	LOCATION
Minnesota	1017 1026	LOCATION

```
<TERM_EXTRACTION_RESULTS> <LIST_TERM_CANDIDATES>
<TERM_CANDIDATE MM_STATUS='1'>
<ID>1</ID>
<FORM>control</FORM>
<LEMMA>control</LEMMA>
<SYNTACTIC_FEATURES>
<SYNTACTIC_CATEGORY>NN-SYNTACTIC_CATEGORY>
<MORPHOSYNTACTIC_FEATURES>
<HEAD>tert665</HEAD>
<NUMBER>2</NUMBER>
<NUMBER_OF_OCCURRENCES>8</NUMBER_OF_OCCURRENCES>
<LIST_OF_OCCURRENCES>
<OCCURRENCE>
<ID>1</ID>
<NP>B-NP</NP>
<DOC>25</DOC>
<SENTENCE>
<START_POSITION>168</START_POSITION>
<END_POSITION>167</END_POSITION>
</OCCURRENCE>
```

The	The	DT	B-NP	O
etiology	etiologic	JJ	I-NP	O
and	and	CC	I-NP	O
epidemiologic	epidemiologic	JJ	I-NP	O
spectrum	spectrum	NN	I-NP	O
of	of	IN	B-PP	O
bronchiolitis	bronchiolitis	NN	B-NP	O

You see the problem.

Added to the other problems, this makes NLP very difficult for scientists to exploit.

Portions of the preceding slides based on:

Kralinger, Martin, Valencia, Alfonso, and Hirschman, Lynette. Linking genes to diseases: mining, extraction, and retrieval applications for biology. *Genome Biology* 9:2, 2008.

Tutorial: Text Mining for Biomedicine: Techniques & Tools, Slides by Sophia Ananiadou, National Centre for Text Mining

Biomedical Natural Language Processing, by Kevin Bratton, Cohen, Dina Demner-Fushman, Benjamins Publishing.

Nancy Ide, Keith Suderman
Vassar College

James Pustejovsky, Marc Verhagen

Brandeis University

Christopher Cieri, Denise DiPersio, Jonathan Wright
Linguistic Data Consortium (Penn)

Eric Nyberg, Di Wang
Carnegie Mellon University

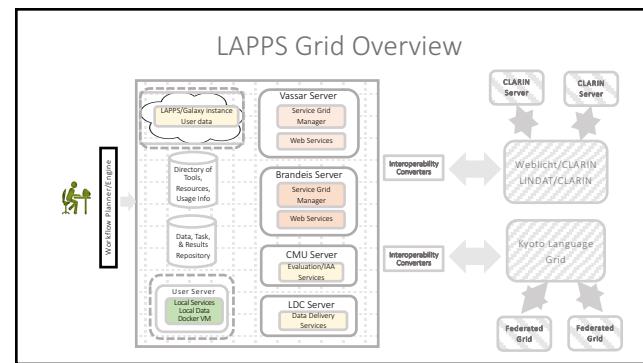
The Language
Applications (LAPPS)
Grid



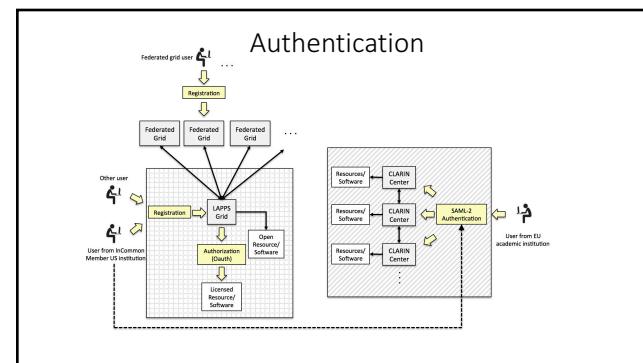
LAPPs/Galaxy Interface

 <http://galaxyproject.org>

- Galaxy is an open, web-based platform designed primarily for **computational biomedical research**
 - **Accessible:** Users without programming experience can easily specify parameters and run tools and workflows
 - **Reproducible:** Galaxy captures information so that any user can repeat and understand a complete computational analysis
 - **Transparent:** Users share and publish analyses via the web and create interactive, web-based documents that describe a complete analysis
- The LAPPs Grid uses the GALAXY framework as a vehicle to combine services of the Language Application Grid
 - Text processing pipelines, components wrapped as services, visualization of component output, evaluation of alternate pipelines, saving and sharing workflows, etc.



The screenshot shows the LAPPs/Galaxy web interface. The main title is "The Language Applications Grid" with the subtitle "An open framework for interoperable NLP web services". On the left, there's a sidebar with various tools and datasets. The main content area displays a search result for "The Language Applications Grid". It includes a list of items such as "1. Fetch documents from language corpora and data from lexicons and other language resources.", "2. Create and apply workflows using tools drawn from several major NLP projects and platforms.", and "5. ZONING ZOOM: Access hundreds of tools and resources available from the Language Grid and other federated grids in Asia, as well as EU CLARIN Nodes/Cards". Below this is a note about license agreements and terms of use.



LAPPs/GALAXY

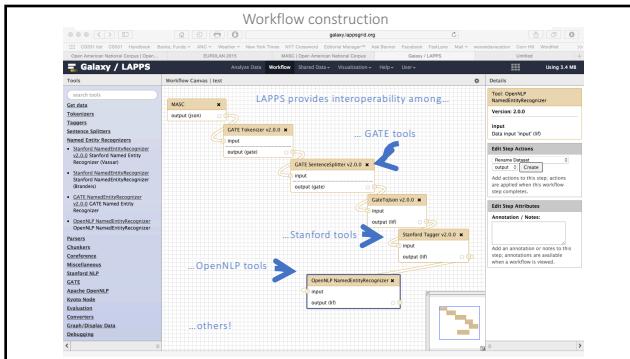
Multiple options for running the LAPPs/Galaxy instance:

1. Use the LAPPs/Galaxy web interface
 - <http://galaxy.lappsggrid.org>
2. Create a local Galaxy instance including:
 - All of Galaxy, or
 - The Galaxy "NLP Flavor" with only LAPPs tools
3. Create a docker image that is a self-contained vm running LAPPs/Galaxy
 - Useful when privacy required, no network connection available, etc.
4. Create a Galaxy instance in the cloud
 - Useful for large datasets, computationally intense applications
 - Useful for tutorials!
 - **HERE: <https://jetstream.lappsggrid.org>**

Interoperability

- The Galaxy bio tools are not interoperable
 - Multiple converters among genomic formats where it is possible to provide them
- LAPPs Grid services are interoperable!

The screenshot shows the Galaxy bio tools interface. It features a sidebar with "Format Converters" and a main panel titled "Galaxy Service" for "TUCF Genomics". The main panel displays a "Galaxy Service" section with a "Generate Sequences" tool and a "Data Services" section listing various datasets and their details.



Brief digression: What is interoperability?

INTEROPERABILITY DEFINED

in·ter·op·er·a·bil·i·ty

\in-tər-ä-p(ə-)rə-'bi-lə-tē\

noun

ability of a system (as a computer system or software) to exchange and make use of information.

INTEROPERABILITY DEFINED

- ▶ **Interoperability** is a characteristic of a product or system, whose interfaces are completely understood, to work with other products or systems, present or future, in either implementation or access, without any restrictions."

Wikipedia

More precisely...

- ▶ Data is in a format that is immediately readable and processable by a given piece of software
- ▶ Output of a given piece of software can be read into and operated upon by another piece of software without modification

Two levels of interoperability

Syntactic interoperability

- Relies on specified data formats, communication protocols, and the like to ensure communication and data exchange

Semantic interoperability

- The content of the information exchange requests are unambiguously defined: what is sent is the same as what is understood

Syntactic interoperability

- ▶ Systems involved can process the exchanged information, but **no guarantee that the interpretation is the same**

▶ EXAMPLE: Data in XML format is readable by software designed to recognize XML



Semantic interoperability

- ▶ Two systems have the ability to automatically interpret exchanged information meaningfully and accurately in order to produce useful results via reference to a common information exchange reference model



Interoperability for language data

- ▶ Syntactic interoperability achieved with common physical formats
 - ▶ Many options: One sentence per line, part-of-speech tag appended to word, XML...

- ▶ Semantic interoperability achieved with common definitions for labeled data
 - ▶ E.g., labels like *noun*, *person*, *date* mean the same to both systems

Not easy! There can be subtle differences of opinion (e.g., should "in the future" be labeled as a DATE? Is "the White House" a LOCATION or an ORGANIZATION in a phrase like "The White House said today...")?
Let alone that people do not agree on the exact definition of noun...

Obstacles

- ▶ Difficult to identify a single representation format that accommodates all kinds of language data and annotations
- ▶ Difficult to get the community to agree, adopt a single standard
- ▶ Need to accommodate legacy data and tools using other formats

Current solution

- ▶ 30 years of development have led to reasonable convergence of practice
- ▶ Key idea:
 - ▶ Instead of defining a single solution, design a universal “pivot” into and out of which other schemes can be easily mapped
 - ▶ For physical formats, requires that the pivot is a serialization of a common abstract data model (directed acyclic graph)
 - ▶ This model underlies UML, ER diagrams, RDF, JSON and JSON-LD, XML, semantic and other kinds of networks...
 - ▶ For semantics, provide a common structured set of terms to which other schemes can be mapped

Once again, this is non-trivial

So, How Does the LAPPs Grid Enable Interoperability?

- **LAPPs Interchange Format (LIF)**
 - Format that allows web services to exchange detailed information about data and its annotations
 - **Syntactic interoperability**
 - handled by **JSON-LD**
 - enforced by the **LIF JSON schema**
 - **Semantic interoperability**
 - enhanced by using the Linked Data aspect of JSON-LD to link to the **LAPPs Web Services Exchange Vocabulary**

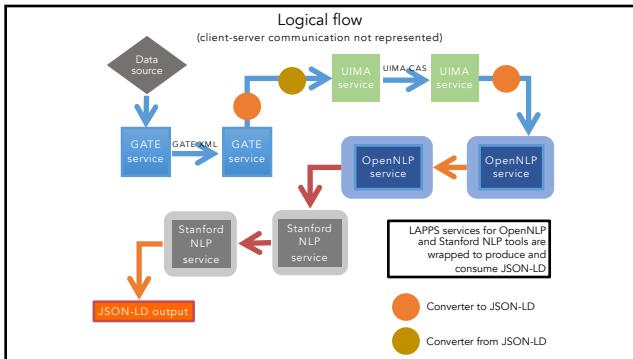
Web Service Communication in LAPPs

1

- Each service in the LAPPs Grid publishes metadata:
- a **discriminator** (**type**) : tells how to interpret the payload
 - a **payload** (typically a utf-8 string)

2

- LAPPs uses JSON-LD as its standard format for the payload
- **Converters** to and from JSON-LD for services that deliver in other formats
 - Some LAPPs services are wrapped to produce and consume JSON-LD



LAPPS Grid Web Service Exchange Vocabulary

- No accepted standard for module description or input/output interchange in the language application domain
- LAPPS Web Service Exchange Vocabulary (WS-EV)**
 - specifies a terminology for a core of linguistic objects and features exchanged among NLP tools that consume and produce linguistically annotated data
 - addresses a need within the community to identify a standard terminology and indicate the relations among them

Design Principles

01 Orthogonal design • Only one entry per concept	02 Lightweight • Easy to find on the web and reference	03 Flexible • Use what you need, add what you need	04 (Arbitrary) decisions about what goes where • Up to this for exchange only • Not confined to the WS-EV terminology or organization internally
--	---	---	--

User definition of objects and features

Options	Many web services will require definition of objects and properties not included in the WS-EV or elsewhere
Options	Provide a URI where a new term or other documentation is defined Add a definition to the WS-EV • Service providers use the name space automatically assigned to them at the time of registration • Avoids name clashes • Makes a distinction between general categories used across services and more idiosyncratic categories

Implementation

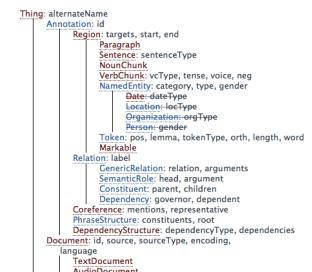
• Bottom-up approach

- Define objects and properties as needed to accommodate LAPPs services as they are added to the LAPPs Grid
- Avoids *a priori* development of a comprehensive standard linguistic type system
- “Minimalist” strategy to provide a simple core set of objects and features
- User capacity to add/replace objects and properties to allow for dynamic typing

LAPPs WS-EV Repository

- <http://vocab.lappsgrid.org>
- Shallow hierarchy of elements
 - Inheritance

LAPPs Exchange Vocabulary Type Hierarchy



Spec for Token

Thing > Annotation > Span > Token

Definition: A string of one or more characters that is a part of a larger sequence of characters. It is the smallest unit of morpho-syntactic labeling apart from dependency annotations.

Annotations: The entire set of tag set used by the part of speech tagger.

Metadata from Annotation:

Properties	Type	Description
Annotation	Annotation or URL	The annotation that produced the annotations.
Label	Label or URL	The documentation of any file or note that were used to identify the annotations.

Properties:

Properties	Type	Description
content	String or URL	The content of the tokens as annotated with the tokens.
contentOffset	String or URL	The offset of the tokens relative to the primary data. (e.g. 0-1000 for a sentence).
contentType	String or URL	Sub-type of the token as word, punctuation, abbreviation, number, article, etc. Similar to a URL referencing a pre-defined descriptor.
end	Integer	The ending offset of the token in the primary data.
length	Integer	The length of the token.

Properties from Span:

Properties	Type	Description
id	Text	ID value of the annotations that make up the token in the primary data.
label	Text	A unique identifier associated with the annotation.
start	Text	The ending offset of-based in the primary data.

Properties from Annotation:

Properties	Type	Description
label	Text	A unique identifier associated with the annotation.

Properties from Thing:

Properties	Type	Description
label	Text	An identifier for the term.

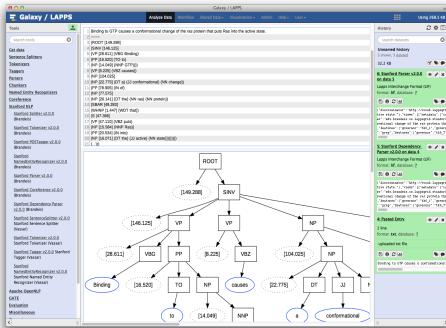
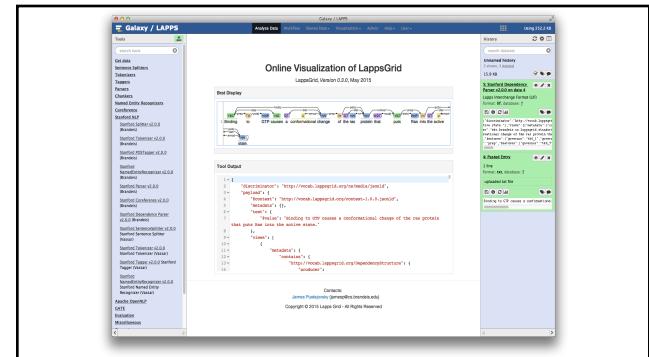
JSON-LD and WS-EV

- References in JSON-LD representation point to URLs providing **definitions** for specific linguistic categories in the WS-EV
- Also point to **documentation** for processing software and rules for processes such as tokenization, entity recognition, etc.
 - Often left unspecified in annotated resources
 - Not required for web service exchange in the LAPPs Grid
 - **BUT** inclusion of such references can contribute to better replication and evaluation of results in the field
 - **Promote best practice!**

JSON-LD and the LAPPs Exchange Vocabulary

```

{
  "@context": "http://vocab.lappsgrid.org/ Base URL for the LAPPs Exchange Vocabulary
  "metadata": {
    "text": "Some of the strongest critics of our welfare system -"
  },
  "views": [
    {
      "metadata": {
        "contains": [
          {
            "tokens": [
              {
                "producer": "org.anc.lapps.stanford.galloway14", Internal LAPPs type defined at http://vocab.lappsgrid.org/LF
                "type": "tokenization:stanford"
              }
            ]
          }
        ],
        "annotations": [
          {
            "type": "Token", Defined at http://vocab.lappsgrid.org/token
            "id": "token",
            "start": 18, Features defined at http://vocab.lappsgrid.org/Token#feature-name
            "end": 22,
            "features": [
              {
                "string": "Some"
              }
            ]
          }
        ]
      }
    }
  ]
}
  
```



EVALUATION IN LAPPs

- CMU has implemented services for state-of-the-art **Open Advancement** techniques
- Provides an unprecedented tool for NLP development**
 - Could take the field to a new level of productivity
- Enables rapid identification of
 - frequent error categories within modules and documents
 - which module(s) and error type(s) have the greatest impact on overall performance
- Used in the development of IBM's Watson to achieve steady performance gains over the four years of its development

Open Advancement in a Nutshell

01

Evaluate multiple possible solutions (tool configurations) for a given problem

- Determine the optimal solution available using given components, resources, and evaluation data

02

Output of the optimal solution subjected to error analysis

- Identify the most frequent errors with the highest impact
- Consider possible enhancements
- Aim to achieve the largest possible reduction in error rate by addressing the most frequent error types

03

Evaluate performance of new configurations

- Determine if a significant improvement has been achieved in comparison with prior baselines

CMU OAQA

- Open advancement for Question Answering
- Analyzes results in/from alternative pipelines

Tool	Reference Outputs	Predicted Outputs	Start	End	Features	Start	End	Features
Get data	362	362	None	None	None	None	None	None
Tokenizers	362	362	None	None	None	None	None	None
Pos Tagger	362	362	None	None	None	None	None	None
Sentence Splitters	362	362	None	None	None	None	None	None
Named Entity Recognizers	362	362	None	None	None	None	None	None
Patterns	None	None	None	None	None	None	None	None
Corference	13	14	IN	IN	NN	13	14	NN
Miscellaneous	15	15	NN	NN	NN	15	15	NN
Stanford CoreNLP	64	65	NN	NN	NN	64	65	NN
GATE	66	67	NN	NN	NN	66	67	NN
Apache OpenNLP	69	70	NN	NN	NN	69	70	NN
OpenNLP Models	74	75	NN	NN	NN	74	75	NN
Evaluation	78	82	CD	CD	NN	78	82	CD
Converters	83	83	NN	NN	NN	83	83	NN
Graph-Dotless Data	94	94	NN	NN	NN	94	94	NN
Workflow	92	97	CD	CD	NN	92	97	NN
• All workflows	98	99	LRB	LRB	NN	98	99	LRB
	102	103	RBB	RBB	NN	102	103	RBB
	104	108	IN	IN	NN	104	108	IN
	108	109	NNP	NNP	NN	108	109	NNP
	132	139	NNP	NNP	NN	132	139	NNP
	139	140	NNP	NNP	NN	139	140	NNP
	141	151	NNP	NNP	NN	141	151	NNP

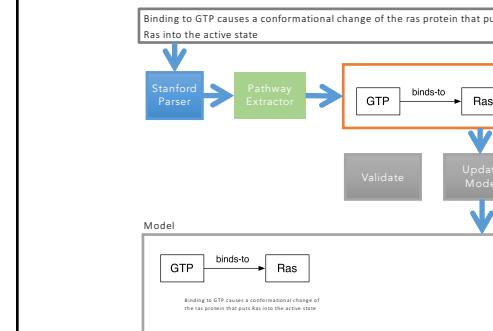
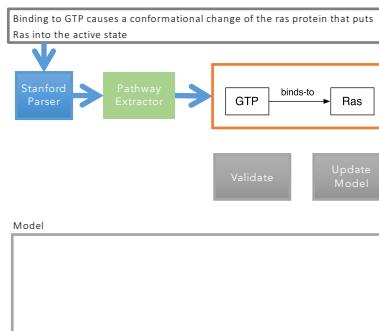
Reference Outputs	Predicted Outputs
Start	Start
End	End
Features	Features
NN	NN

BENEFITS OF THE LAPPS/GALAXY COLLABORATION

- Galaxy contains a huge number of tools for analyzing genomic and other biomedical data
- LAPPS includes tools to perform NLP analyses
- Combining LAPPS services with Galaxy tools can provide data mined from the vast stores of biomedical literature (Biomed, PubMed, PLOS, etc.)
 - interaction between model and observations (pathway steps), can be mined from new text
- BIONLP meets bio-analysis!

Example

Binding to GTP causes a conformational change of the ras protein that puts Ras into the active state. GTP-bound ras binds to the raf protein kinase. This binding of raf to ras has the effect of activating the raf kinase and localizing the raf kinase to the cell membrane. Activated raf now phosphorylates and activates the Mek1 kinase. The Mek1 kinase then phosphorylates the ERK kinase on both threonine and tyrosine residues which activate ERK kinase activity. The phosphorylated ERK protein then translocates to the nucleus where it regulates gene expression in part by phosphorylating the Elk1 transcription factor. Phospho-Elk then upregulates the gene expression of target genes such as the proto-oncogene c-fos. The entire signaling cascade is terminated by the intrinsic GTPase activity of ras which hydrolyzes the bound GTP into GDP, thus returning ras to the GDP bound state where it releases bound raf. The GTPase activity of ras is accelerated by interaction with another protein called GAP. The oncogenic rasv12 mutant has diminished GTPase activity and therefore stays in the active GTP bound state constitutively. Deletion of GAP or the related NF1 genes will also enhance ras activity by slowing the rate of ras-GTP hydrolysis.





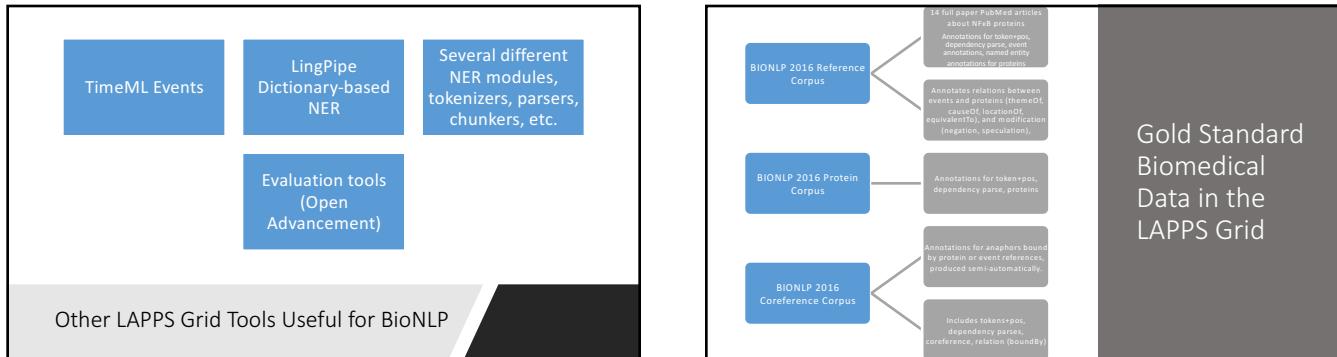
Current Activities

- New NSF grant
 - Collaboration between Vassar College and Penn State University (and other Galaxy Principal Investigators) to
 - Develop a robust and interoperable set of tools, ready-made workflows, etc. for mining biomedical publications
 - Rely on researchers to test and provide feedback as we develop
 - Provide seamless integration of text mining capabilities and the vast array of tools provided in Galaxy
- Collaboration with the US government Centers for Disease Control and Food and Drug Administration to adapt the LAPPs Grid for summarization and mining of clinical reports

Penn BioTokenizer
Biomedical NER

Gene annotator

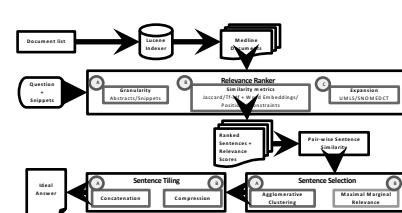
BioNLP-oriented Tools in the LAPPs Grid



BioASQ

- Large-scale biomedical semantic indexing and question answering challenge, started in 2013
- Phase B of Task 5b of the BioASQ challenge: Biomedical Q/A
 - Involved answering the following types of questions
 - Factoid
 - List
 - Yes/ No
 - Summary
 - Exact answers for Factoid, List, Yes/ No
 - Ideal answers for all questions types
 - Five test sets, comprising of 100 questions each

Pipeline Architecture of Winning System



Hands-on Exercises

01

- Using the LAPPs Grid
- Learn how to import data, execute workflows, visualize results, etc.
 - Explore various tools, including parsers, named entity recognizers

02

- Exploring biomedical data and tools
- Import data without annotations, annotate with LAPPs Grid
 - Tokenizer, protein annotator, event annotator
 - Dictionary-based NER

03

- Evaluating automatically-produced results
- Use LAPPs Grid evaluation services to compare gold standard annotations with automatically-generated ones