

The Language Application Grid and Galaxy

Nancy Ide, Keith Suderman Vassar College

James Pustejovsky, Marc Verhagen Brandeis University

Christopher Cieri Linguistic Data Consortium

Eric Nyberg Carnegie-Mellon University

The LAPPS Grid

- A framework to
 - enable language service discovery, composition, and reuse
 - For both NLP researchers and others (who may use predeveloped composite services)
 - promote sustainability, manageability, usability, and interoperability of NLP components
- Based on the service-oriented architecture (SOA)
 - Web-oriented version of the "pipeline" architecture for sequencing loosely-coupled linguistic analyses

The LAPPS Grid

- Provides interoperable access to
 - Wide array of NLP processing tools and components
 - language resources such as mono- and multi-lingual corpora and lexicons
- Enables pipelining tools to create custom NLP applications and "black box" composite services
- Provides an open advancement (OA) framework for component- and application-based evaluation
- Will enable easy navigation through licensing issues
- Actively pursuing creation of an interoperable global network of grids and frameworks

LAPPS Grid Overview



Interoperability

• LAPPS Interchange Format (LIF)

- allows services to exchange information
- 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
 - Not yet-another-repository! Linked to others where possible

LAPPS Exchange Vocabulary Type Hierarchy



Home

Thing > Annotation > Region > Token

 Definition
 A string of one or more characters that serves as an indivisible unit for the purposes of morpho-syntactic labeling (part of speech tagging).

 Similar to
 http://www.isocat.org/datcat/DC-1403

 URI
 http://vocab.lappsgrid.org/Token

Metadata

Properties	Туре	Description
posTagSet	String or URI	The definition of the tag set used by the part-of-speech tagger.

Metadata from Annotation

Properties	Туре	Description
producer	List of URI	The software that produced the annotations.
rules	List of URI	The documentation (if any) for the rules that were used to identify the annotations.

Properties

Properties	Туре	Description
pos	String or URI	Part-of-speech tag associated with the token.
lemma	String or URI	The root (base) form associated with the token. URI may point to a lexicon entry.
tokenType	String or URI	Sub-type such as word, punctuation, abbreviation, number, symbol, etc. Ideally a URI referencing a pre-defined descriptor.
orth	String or URI	Orthographic properties of the token such as LowerCase, UpperCase, UpperInitial, etc. Ideally a URI referencing a pre-defined descriptor.
length	Integer	The length of the token
word	String	The surface string in the primary data covered by this Token.

Properties from Region



All tools' input/output formats mapped into and out of LIF

Linguistic categories etc. mapped to WSEV



LAPPS Web Service Exchange Vocabulary

- Specifies a terminology for a core of linguistic objects and features exchanged among NLP tools that consume and produce linguistically annotated data
- Linked wherever possible to existing repositories such as ISOCat (CLARIN Concept Repository), schema.org, FoLiA categories, etc.
- References in JSON-LD representation point to URIs providing definitions for specific linguistic categories in the WS-EV

Current collaborations/projects

• Federation of Service Grids

- LAPPS Grid, Language Grid (Kyoto University, Japan), NECTEC (Thailand), University of Indonesia, Xinjiang University (China), ELRA Grid
- Access to all tools, applications, and resources on any grid through any portal

LAPPS/CLARIN

- CLARIN/WebLicht (Tubingen) and LINDAT/CLARIN (Prague)
- Mellon Foundation proposal to create a trust network between LAPPS and CLARIN

OpenMinted

- Advisory board—work together on harmonization
- LAPPS Grid used in
 - DARPA LORELEI project for under-resourced languages
 - HathiTrust Research Center (HTRC) text mining project
 - Multi-day course for government analysts
 - Undergraduate and graduate CL courses at CMU, Brandeis, Vassar
 - ? IBM Watson

LAPPS Galaxy Interface



http://galaxyproject.org

- LAPPS recently adopted the GALAXY workflow engine as a front end for construction of pipelines etc.
 - open, web-based platform developed for computational genomics/biomedical research

Why Galaxy?

- Accessible: Accommodates users with a broad range of expertise (non-computational to expert programmer)
- Reproducible: Galaxy captures information so that any user can understand and repeat a complete computational analysis
- **Transparent:** Users share and publish analyses via the web and create interactive, web-based documents that describe a complete analysis
- Well-developed, supported, open !

Galaxy

- Available as
 - A free public web server supported by the Galaxy Project (Johns Hopkins, Penn State)
 - Includes widely used bioinformatics tools
 - Users save histories, workflows, and datasets on the server, all can be shared with others
 - Open source software that can be downloaded and installed locally or in a cloud, and customized to address specific usages
 - Public web servers hosted by other organizations

LAPPS/Galaxy

- Multiple options for running a LAPPS/Galaxy instance
 - Use the LAPPS/Galaxy web interface
 - http://galaxy.lappsgrid.org
 - Create a local Galaxy instance:
 - Clone our fork of the Galaxy project, or
 - Run the LAPPS Grid appliance
 - a series of **docker images** that is a self-contained vm running Galaxy and all LAPPS services
 - useful when privacy required, no network connection available, etc.
 - Create an instance in the cloud
 - Useful for large datasets, computationally intense applications

Replicability and Sharing

- The field of NLP research and development has been plagued by a chronic lack of replicability of results
 - A great deal of re-inventing of the wheel and wasted effort
 - Evaluation of results hampered when details of a study (including versions and parameters for data, software) are not included in papers
- The field of NLP is still hampered by a lack of widespread sharing of resources that are the basis of research results

Galaxy as Promoter of Open, Replicable Research

- The field of NLP research and development has been plagued by a chronic lack of replicability of results and information about provenance
- Galaxy provides for:
 - automatic recording of inputs, tools, parameters and settings used for each step in an analysis in a publicly viewable history
 - sharing datasets, histories, and workflows via web links
 - creation of custom web-based documents to communicate about an experiment or result
- Encourages open publication, sharing of data and results
- Fosters replicability and reuse by providing a physical infrastructure to enable it

LAPPS/Galaxy

- LAPPS/Galaxy components are LAPPS web services
- Access to 100+ LAPPS services plus those of federated partners
- Interoperability
 - LAPPS services among each other
 - LAPPS services and Galaxy components
 - handled by LAPPS converters

Galaxy / LAPPS

1

0

Using 0 bytes

C 🗘

0

Tools

Local Data

MASC

<u>Gigaword</u>

Tokenizers

Sentence Splitters

Taggers

Named Entity Recognizers

Parsers

Chunkers

Stanford NLP

GATE

Apache OpenNLP

Lingpipe

DKPro Core

DBpedia

Evaluation

Manual Conversion

Miscellaneous

Graph/Display Data

Workflows

All workflows

LAPPS Grid

A Framework for Rapid Adaption and Reuse.

Work In Progress

Analyze Data

Many of the services here are undergoing active development and their behaviour is likely to change without notice.

Workflow Shared Data - Visualization - Help - User -

History

0 b

Unnamed history

1 This history is empty. You can

from an external source

load your own data or get data

Welcome to the LAPPS Grid Galaxy instance. Through this Galaxy instance you can:

- 1. Fetch documents from the MASC or Gigaword corpora.
- 2. Create processing pipelines with tools from:
 - GATE
 - 2. Apache OpenNLP
 - 3. Stanford NLP

Simple Tutorial

If you have a good understanding of how Galaxy works you can run the following tools in order:

- 1. Get data -> MASC
- 2. From the GATE menu ->
 - 1. Tokenizer
 - 2. Sentence Splitter
 - 3. Part of speech tagger
- 3. From the History panel select ->
 - 1. Edit attributes
 - 2. Convert Format (there is only one converter, so just run it)
- 4. Tools -> Word Count
- E Expand the output calest the Visualiza icon and then Charte

Workflow construction

		galaxy.lappsgrid.org	
CS331 list CS331 Handbook B	anks, Funds 🖌 ANC 🗸 Weather 🖌 New York Times 🕅	NYT Crossword Editorial Manager™ Ask Banner Facebook FastLan	e Mail ~ weneedavacation Corn Hill WordNet >
Open American National Corpus Open	EUROLAN 2015 MA	ASC Open American National Corpus Galaxy / LAPPS	Untitled +
🚍 Galaxy / LAPPS	Analyze Data Workflo	w Shared Data - Visualization - Help - User -	Using 3.4 MB
Tools	Workflow Canvas test		Details
search tools	LAPP:	S provides interoperability	Tool: OpenNLP NamedEntityRecognizer
<u>Get data</u>	MASC × amono	J	Version: 2.0.0
<u>Tokenizers</u>	output (json)		
Sentence Splitters	GATE Tokenizer v2.0.0 ×	GATE	input Data input 'input' (lif)
Named Entity Recognizers Stanford NamedEntityRecognizer v2.0.0 Stanford Named Entity	output (gate)	tools	Edit Step Actions
Recognizer (Vassar)		GATE SentenceSplitter v2.0.0 ×	Rename Dataset ≎ output ≎ Create
 <u>Stanford NamedEntityRecognizer</u> Stanford NamedEntityRecognizer (Brandeis) 		output (gate)	Add actions to this step; actions are applied when this workflow
<u>GATE NamedEntityRecognizer</u> <u>v2.0.0</u> GATE Named Entity Recognizer		GateToJson v2.0.0 ×	Edit Step Attributes
OpenNLP NamedEntityRecognizer OpenNLP NamedEntityRecognizer		Stanford tool	Annotation / Notes:
Parsers		Stanford Tagger vi	.0.0 🗙
<u>Chunkers</u>		input	Add an annotation or notes to this
<u>Coreference</u> <u>Miscellaneous</u>	OpenNLP 🔍	output (lif)	step; annotations are available when a workflow is viewed.
Stanford NLP	tools		
GATE		OpenNLP NamedEntityRecognizer 🗙	
Apache OpenNLP		input	
Kyoto Node		output (lif)	
Evaluation			
Cranh/Display Data	others!		
Debugging			
<	H		

Step-by-step analysis

000	Galaxy / LAPPS	R _M
🗧 Galaxy / LAPPS	Analyze Data Workflow Shared Data - Visualization - Admin Help - User -	Using 246.9 KB
Tools	Download data directly from web or unload files from your disk	History 📿 🌣 🗔
(search tools		search datasets
<u>Get data</u> <u>Upload File</u> from your computer	You added 1 file(s) to the queue. Add more files or click 'Start' to proceed.	Unnamed history 3 <u>deleted</u>
Gigaword (LDC) Retrieve a	Name Size Type Genome Settings Status	10.6 KB
MASC Retrieves a single document from MASC.	New File 0.1 KB Auto-det Q unspecified (?) The second sec	This history is empty. You can load your own data or get data from an external source
List MASC Lists all keys in MASC.	Binding to GTP causes a conformational change of the ras protein that puts Ras into the active state.	
MASC Load Retrieves a single of document from MASC.		
Sentence Splitters	Sir	
Taggers		
Parsers	if y	
Chunkers		
Named Entity Recognizers		
Coreference	Type (set all): Auto-datact Q Genome (set all): Uncoarified (2)	
Anache OpenNI P		
GATE		
Evaluation	Choose local file Paste/Fetch data Start Pause Reset Close	
Miscellaneous		
Converters	7. Click the braw button to marvel at the simple word count histogram ;-)	
Graph/Display Data		
Debugging		
Workflows	Caveats	
<u>All workflows</u>	While Galaxy ansures that the inputs to a processing tool use the correct file format (CATE tools require CATE/XM), the other tools require LIE) it does not	
	currently ensure the document contains the required annotation types. For example; part of speech taggers typically require the input document to contain	
	Token annotations.	
	Galaxy is an open, web-based platform for data intensive biomedical research. The Galaxy team is a part of BX at Penn State, and the Biology and Mathematics and	
	Computer Science departments at Emory University. The Galaxy Project is supported in part by NHGRI, NSF, The Huck Institutes of the Life Sciences, The Institute for CyberScience at Penn State, and Emory University.	
2		

000	Galaxy / LAPPS		R <u>w</u>
💶 Galaxy / LAPPS	Analyze Data Workflow Shared Data - Visualization - Admin Help - User -	Us	ing 252.2 KB
Tools		History	≈ ≎ 🗆
search tools	E Stanford Dependency Person v3.0.0 on date 4	search datasets	8
<u>Get data</u>	5: stantord Dependency Parser V2.0.0 on data 4	Unnamed history	
<u>Sentence Splitters</u> Tokenizers	from 'running' to 'finished' if completed successfully or 'error' if problems were encountered.	2 snown, 3 <u>deleted</u> 15.9 KB	
Taggers Parsers		5: Stanford Dependency Parser v2.0.0 on data 4	• # ×
<u>Chunkers</u>		4: Pasted Entry	• # X
Named Entity Recognizers Coreference Stanford NI P		1 line format: txt , database: <u>?</u>	
Stanford Splitter v2.0.0		uploaded txt file	
(Brandeis) Stanford Tokenizer v2.0.0		8821	•
(Brandeis)		Binding to GTP causes a c	onformational
<u>Stanford POSTagger v2.0.0</u> (Brandeis)			
<u>Stanford</u> <u>NamedEntityRecognizer v2.0.0</u> (Brandeis)			
<u>Stanford Parser v2.0.0</u> (Brandeis)			
<u>Stanford Coreference v2.0.0</u> (Brandeis)			
<u>Stanford Dependency Parser</u> <u>v2.0.0</u> (Brandeis)			
<u>Stanford SentenceSplitter v2.0.0</u> Stanford Sentence Splitter (Vassar)			
<u>Stanford Tokenizer v2.0.0</u> Stanford Tokenizer (Vassar)			
<u>Stanford Tagger v2.0.0</u> Stanford Tagger (Vassar)			
<u>Stanford</u> <u>NamedEntityRecognizer v2.0.0</u> Stanford Named Entity Recognizer (Vassar)			
Apache OpenNLP			
GATE Evaluation			
Miscellaneous			
<			>



00	Galaxy / LAPPS		E _M
🚍 Galaxy / LAPPS	Analyze Data Workflow Shared Data - Visualization - Admin Help - User -	III (lsing 268.5 KB
Tools	1 job has been successfully added to the gueue - resulting in the following datasets:	History	€ ♥ 🗆
search tools	6: Stanford Parser v2.0.0 on data 5	search datasets	0
<u>Get data</u>	You can check the status of queued into and view the resulting data by refreshing the History game. When the int has been run the status will change	Unnamed history	
Sentence Splitters	from 'running' to 'finished' if completed successfully or 'error' if problems were encountered.	32.2 KB	
Taggers		C: Stanford Davage 20	
Parsers		on data 5	<u>,</u>
Chunkers		5: Stanford Dependency	
Named Entity Recognizers		Parser v2.0.0 on data 4	
Stanford NLP		Lapps Interchange Forma format: lif. database: ?	ıt (LIF)
<u>Stanford Splitter v2.0.0</u> (Brandeis)		B 0 2 III	•
<u>Stanford Tokenizer v2.0.0</u> (Brandeis)		<pre>{"discriminator":"http:/ tive state."},"views":[{ er":"edu.brandeis.cs.lap unities1</pre>	/vocab.lappsgri "metadata":{"cc psgrid.stanford
<u>Stanford POSTagger v2.0.0</u> (Brandeis)		<pre>rmational change of the ,"features":{"governor": ":"prep","features":{"go</pre>	ras protein the "tk0_1","govern vernor":"tk0_7"
<u>Stanford</u> NamedEntityRecognizer v2.0.0 (Brandeis)		4: Pasted Entry	• # ×
<u>Stanford Parser v2.0.0</u> (Brandeis)		format: txt , database: <u>?</u>	
<u>Stanford Coreference v2.0.0</u> (Brandeis)			•
<u>Stanford Dependency Parser</u> <u>v2.0.0</u> (Brandeis)		Binding to GTP causes a	conformational
<u>Stanford SentenceSplitter v2.0.0</u> Stanford Sentence Splitter (Vassar)			
<u>Stanford Tokenizer v2.0.0</u> Stanford Tokenizer (Vassar)			
<u>Stanford Tagger v2.0.0</u> Stanford Tagger (Vassar)			
<u>Stanford</u> <u>NamedEntityRecognizer v2.0.0</u> Stanford Named Entity Recognizer (Vassar)			
Apache OpenNLP			
GATE			
Evaluation			
Miscellaneous			
<			>

00	Galaxy / LAPPS		R <u>s</u>
🔁 Galaxy / LAPPS	Analyze Data Workflow Shared Data - Visualization - Admin Help - User -		Using 268.5 KB
Tools	1 Binding to GTP causes a conformational change of the ras protein that puts Ras into the active state.	History	2 * 🗆
search tools	2 ~~~~ 3 (ROOT [149 288]	search datasets	0
<u>Get data</u>	4 (SINV [146.125]	Unnamed history	
Sentence Splitters		3 shown, 3 <u>deleted</u>	
Tokenizers	7 (typ [14.049] (typ GTP))))	32.2 KB	
Taggers	8 (VP [8.225] (VBZ causes))	6: Stanford Parser v2 (
Parsers	9 (NP [104.025]	on data 5	
Chunkers	10 (NP [22,7/5] (D1 a) (JJ conformational) (NN change))	Lapps Interchange Forn	nat (LIF)
Named Entity Recognizers	2 (NP (7.575) 2 (NP (7.575)	format: lif, database: ?	
Coreference	13 (NP [26.141] (DT the) (NN ras) (NN protein))	BA2	
Stanford NLP	14 (SBAR [49.283]		
Stanford Splitter v2.0.0	15 (WHNP [1.447] (WDT that))	<pre>{"discriminator":"http: tive state."}."views":'</pre>	://vocab.lappsgri [{"metadata":{"cc
(Brandeis)	10 (5 (47.300) 17 (VP (47.110) (VBZ puts)	er":"edu.brandeis.cs.ls	appsgrid.stanford
Stanford Tokenizer v2.0.0	18 (NP [15.584] (NNP Ras))	rmational change of the "features":{"governor"	e ras protein the ':"tk0 1","govern
(Brandeis)	19 (PP [20.534] (IN into)	":"prep","features":{"c	governor":"tk0_7'
Stanford POSTagger v2.0.0	20 (NP [16.071] (DT the) (JJ active) (NN state)))))))) 21 (, .)))		
(Brandeis)		<u>S: Stanford Dependent</u> Parser v2.0.0 on data	<u>cy</u> (*) * *
Stanford	2007	Lapps Interchange Forn	nat (LIF)
NamedEntityRecognizer v2.0.0 (Brandeis)		format: lif, database: ?	
Stanford Parser v2.0.0		B 0 2 m	۰ ا
(Brandeis)		{"discriminator":"http:	://vocab.lappsgri
Stanford Coreference v2.0.0	149 2881 SINV	tive state."},"views":[er":"edu.brandeis.cs.ls	[{"metadata":{"cc appsgrid.stanford
(Brandeis)		rmational change of the	e ras protein the
Stanford Dependency Parser		,"features":{"governor" ":"prep","features":{"ç	':"tk0_1","govern governor":"tk0_7'
<u>v2.0.0</u> (Brandeis)			
Stanford SentenceSplitter v2.0.0		4: Pasted Entry	• / ×
(Vassar)		1 line	
Stanford Tokonizer v2.0.0		format: txt , database: <u>?</u>	2
Stanford Tokenizer (Vassar)		uploaded txt file	
<u>Stanford Tagger v2.0.0</u> Stanford Tagger (Vassar)	[28.611] VBG PP [8.225] VBZ [104.025] NP	802Ш	•
Stanford NamedEntityRecognizer v2.0.0		Binding to GTP causes	a conformational
Stanford Named Entity			
Recognizer (Vassar)	Binding [16 520] TO NP causes [22 775] DT 11 N		
Apache OpenNLP			
GATE			
Evaluation			
Miscellaneous			
	to [14.049] NNP (a) conformational		
			>

Evaluation in LAPPS/Galaxy

- CMU has implemented services for state-of-the-art
 Open Advancement techniques
- 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

Analyzes results in/from alternative pipelines



- Can be comparison to gold standard, or comparison to another pipeline or pipelines
- Potentially any number of pipelines can be compared
 - CMU working on methods for finding an optimal solution among all multiple possible paths

🚍 Galaxy / LAPPS



MASC gold standard vs. Stanford NEs

Precision: 0.391 Recall: 0.473 F1: 0.428

Reference Outputs			Predicted Outputs			
Start	End	Features	Start	End	Features	Text
70	80	LOCATION				Asia Minor
			70	74	LOCATION	Asia
			82	85	DATE	now
93	99	LOCATION	93	99	LOCATION	Turkey
104	110	LOCATION	104	110	LOCATION	Greece
145	151	LOCATION	145	151	LOCATION	Aegean
389	402	LOCATION				Mediterranean
429	438	DATE				7000 b.c.
			429	433	DATE	7000
			434	438	DATE	b.c.
			478	481	DATE	now
490	494	LOCATION	490	494	LOCATION	Iran
647	657	LOCATION				Aegean Sea
			647	653	LOCATION	Aegean
			654	657	LOCATION	Sea
738	743	LOCATION	738	743	PERSON	Milos
899	904	LOCATION	899	904	PERSON	Milos

Potential benefits of 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 allow for analysis of data mined from the vast stores of biomedical literature (Biomed, PubMed, PLOS, etc.)
- 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 GTP, 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.





Model









LAPPS/GALAXY Demo

P31 LR Infrastructures and Architectures Thursday May 26 14:55-16:35 Poster Area 1

