# Data search and discovery in language data spaces: challenges and solutions

Stelios Piperidis, ILSP/Athena RC spip@athenarc.gr

http://www.european-language-grid.eu, http://www.ilsp.gr 03-06-2022 Data Spaces and Semantic Interoperability





#### Institute for Language and Speech Processing

- Information Technology in Greece
- Single dedicated Language and Speech Institute in Greece

#### Ideas and insights at this meeting from the viewpoint of

- Data "producer" for specific language technology development/applications, e.g. machine translation, aspect-based sentiment analysis
- LT developer e.g. text analytics and information extraction, machine translation, speech technologies (recognition/synthesis/analytics), sign technologies
- Platform and infrastructure design, development and operation



• Part of Athena R.C., the biggest research center dedicated to Computer Science and



# A bit of data spaces/infrastructures/platforms related history

#### META SHARE

LEARN - DISCOVER - PARTICIPATE - CONNECT - LOGIN



CLARIN Virt         Welcome to the VLO!         Use the search bar below to start sorwse everything and use facets         See all records	ual Language Observatory searching through hundreds of thousands of language resources, or continue to to narrow down to your area of interest or discover new resources.	
Search through 1,030,321 records		Q

Virtual Language Observatory Search Contributors Help









$\leftarrow$	$\rightarrow$	G	🔒 li	ve.european-language-grid.eu/catalogue/search/italian%20newspaper?&elg_integrated_services_and_data_term=ELG%20host	ed%20data	ର୍	B	☆	Ľ
	14	E	EUROP LANGU GRID RELEASE 3	EAN AGE	Catalogue 뇌	Documentation &	Media	a V	Ał
	itali	an nev	wspaper			Search		?	

Language resources & technologies	~	-
Languages	~	9
Media types	~	0
Licences	~	ŀ
Conditions of use	~	
ELG integrated services and data	^	
– ELG hosted data (51)		5

Clear all filters 🛞

51 search results for italian newspaper	
ELG hosted data 😣	
Version: 1.0.0	39 views
	66 downloads
The CHANGE-IT dataset contains approximately 152,000 article-headline pairs,	hosted in ELG
collected from two Italian newspapers situated at opposite ends of the political	
spectrum, namely la Repubblica (left) and Il Giornale (right), 🗸	
Keywords: style · natural language generation · news · newspaper headlines	
Language: Italian	
Licence: Creative Commons Attribution Non Commercial Share Alike 4.0 International	
MLCC Multilingual and Parallel Corpora	20 views
	0 downloads
The MLCC text corpus has two main components - one set to allow comparable	hosted in ELG
studies to be carried out in different languages and one set as the basis for	
translation studies. The first set is referred as the Polylingual D $\checkmark$	



1 😩	ELG - CHANGE-IT dataset	age-grid.eu/catalogue/corpus/	7373				Q
© bout \	EUROPEAN LANGUAGE GRID RELEASE 3					Catalogue	e 뇌 Documentation 8
ews ads 3	CHANGE-IT dataset CHANGE-IT Version: 1.0.0 hosted in ELG Keyword style natural language generation news newspaper headlines style transfer	Intended application Natural Language Generation		Corpus subclass		Cite resource De Mattei, Lorenzo; Gatt Cafagna, Michele (2021) (Text corpus)]. De Mattei Nissim, Malvina; Cafagn Cite all versions De Mattei, Lorenzo; Gatt Cafagna, Michele (2021) De Mattei, Lorenzo; Gatt Cafagna, Michele https:/	t, Albert; Dell'Orletta, F . CHANGE-IT dataset. \ i, Lorenzo; Gatt, Albert a, Michele https://doi.o t, Albert; Dell'Orletta, F . CHANGE-IT dataset. [ t, Albert; Dell'Orletta, F //doi.org/10.57771/y2r
	Download						
ews	The CHANGE-IT dataset contains approxiends of the political spectrum, namely la has been used in the context of the CHAN	mately 152,000 article-headline pai Repubblica (left) and Il Giornale (ri <u>c</u> NGE	rs, collected from t ght), with the two r	wo Italian newspapers situated newspapers equally represented.	at opposite The dataset	Share	in 🗈
3	Corpus part					Views	Downloa
	TEXT La	inguage				39	66
		My gric Catalogue 🏼 Docu	EU LA GR REL	ROPEAN NGUAGE ID EASE 3			
	Cita Ge Ser Ulr	all versions rmann, Ulrich (2020, February 28 vice for German to English. [Soft ich https://doi.org/10.57771/qt44	Η πλειονότητα τω «root αντιλι root νεειε πλειονότητα nsubj NOUN Η Ελλήνω nmod	ν Ελλήνων αντιλαμβάνεται ότι το κλίμα αμβάνεται αλλάζει ccomp VERB v ότι κλίμα δραστικά αποτελ mark nsubj advmod conj	του πλανήτη αλλάζει	δραστικά , αποτελώντας ένα μ	εγάλο κίνδυνο για την ανθρ
e wolle "unter o herrscht seit Ta esn't want to co Champions Lea	dem Druck nicht zusammenbrechen": In der französische ägen die öffentliche Debatte. ollapse under the pressure: the French government is doi ague final in Paris has dominated the public debate for d	n Regierung läuft knapp zehn Tage vor ing nothing ays.	DET PROPN TWV det DET	I SCONJ NOUN ADV VEŔB To πλανήτη , det noun punct DET NOUN PUNCT Tou ένα det det DET DET DET	kívõuvo obj NOUN μεγάλο amod ADJ Viα την case det ADP DET	οτητα επιπτώσεις είναι ήδη nsubj NOUN AUX ADV	και θα είναι ακόμα cc aux cop advino CCONJ AUX AUX ADV σ cc cc cc cc cc cc cc cc cc cc cc cc c
				•			
2	<b>.</b>						



4

# **Asset descriptions**

- LT community members
  - can document their assets by providing formal descriptions compliant with a dedicated metadata schema
  - describe & provide access (remote or integrated in ELG) LRTs
  - describe organisations & projects
  - linked with each other
  - must register and get authenticated as providers

		RT	
		H	Tool/Serv
		Н	Datase (Corpus
		H	Lexical/ Conceptu Resourc
		L	Languag Descripti
Overview Download/	Run Try out		Code samples
Keyword machine translation translation multilingual	Intended application Machine Translation		
Language	Function		Language English
Input content resource Language Finnish Processing resource type user input text	Function Function Machine Translation Language dependent true		Language English Processing r output tex
Input content resource Language Finnish Processing resource type user input text Character encoding UTF-8	Function Function Machine Translation Language dependent true		Language English Processing r output tex Character er UTF-8







### ELG in the wider LT and AI ecosystem

- ELG is **building bridges** to existing platforms/infrastructures
  - Mainly in terms of metadata-based descriptions
  - Based on open protocols (OAI-PMH), or **APIs** offered by the platform or infrastructure providers
  - Respecting their own policies
- ELG also as infrastructural arm of ELE
  - using a mixture of automatic and collaborative population of the ELG catalogue













Connecting Europe Facility ELRC-SHAR





















#### **Community repositories**



**HELG** 03-06-2022, Vienna - Data Spaces and Semantic Interoperability



OAI-PMH Client (ELG)

Native metadata repository





### **ML/NLP repositories**



**H**ELG 03-06-2022, Vienna - Data Spaces and Semantic Interoperability





Ingestion

Loading of resulting metadata records to ELG

Conversion

Conversion to ELG metadata

#### **General repositories**





**WELG** 03-06-2022, Vienna - Data Spaces and Semantic Interoperability





#### Ingestion

Loading of resulting metadata records to ELG

#### Validation

Targeted inspection and metadata enrichment

## **Metadata aggregation - interoperability challenges**

- Different types (goals/roles) of repositories: institutional / disciplinary / general repositories
- Different metadata schemas due to different purposes (e.g., preservation, sharing, processing, ...), different resource types, different user needs & practices
- Different platform/infrastructure architectures from fully federated to fully centralised
- Different **policies**, access rights, user access scenarios
- Different data types & sets of values for similar elements (e.g. "language" value as free text vs different **controlled vocabularies**)
- Different granularity level of metadata schema and unit of data (e.g. language data for anthropological studies vs language tech development)
- Different community/application needs for specific elements (e.g. "language" and "resource type")

**H**ELG 03-06-2022, Vienna - Data Spaces and Semantic Interoperability

•



10

### **Elements of solutions?**

- "Open shared semantic space" for vocabularies

  - A common place for vocabularies and their concepts (cf. <u>Fair semantics recommendations</u>) • Linking vocabulary concepts (properties, classes, instances) with semantic relations clearly documented
    - Starting from elements of standard vocabularies (DCAT, schema.org, DC, etc.) linked between them
  - Elements/Values from community-specific vocabularies
    - maintained and curated by communities
- links when possible with other elements (in other domain specific or general vocabularies) • Technically, cross-platform sharing of metadata records
  - cross-platform search instead of a central single endpoint –pros & cons
  - platforms/repos expose their full metadata and each search API selects only those of interest
  - metadata converted (search API exposing catalogue) through the open shared vocabularies

11



#### European Language Grid

#### Thank you!



The European Language Grid has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement № 825627 (ELG). Stelios Piperidis, ILSP/ Athena RC <a href="mailto:spip@athenarc.gr">spip@athenarc.gr</a>

03-06-2022, Data Spaces and Semantic Interoperability http://www.european-language-grid.eu