

Exploring Open Information Extraction for the Portuguese language: An integrated monolithic approach in Cloud environment

Augusto Sampaio Barreto Daniela Barreiro Claro
FORMAS Research Center - Institute of Computing
Federal University of Bahia – Salvador - Bahia - Brazil
augusto.barreto@ufba.br dclaro@ufba.br

Abstract

This work addresses Open Information Extraction (Open IE) as a crucial area for structuring unstructured data, aiming to identify and represent information through triples. The Open IE approach proposes a domain-independent paradigm that extracts potential relationships between entities using generalized patterns. Despite advancements, the lack of studies in the Portuguese language emphasizes the need to explore specific techniques. This is worse regarding easy access to take advantage of triples extracted by OpenIE models. In this context, we present an integrated hub of services in the Open IE domain, allowing users to extract triples and compare their results. This hub is developed under a monolithic architecture with the Django framework, coupled with deployment on Google Cloud to reinforce the efficiency and adaptability of inclusion and removal of services. Within this framework, it would be possible for non-computer experts to use the advantages of OpenIE triples from the Portuguese languages.

1 Introduction

Open Information Extraction (Open IE) is a crucial area that aims to structure data from unstructured sources, with the main goal of identifying and representing information through triples that express relationships. The Open IE approach represents a domain-independent extraction paradigm, using generalized patterns to extract all potential relationships between entities (Etzioni et al., 2008).

Triples are essential for capturing the meaning of information present in unstructured data. Despite advances in the last decade, most of this progress has focused on the English language, with few studies dedicated to Portuguese in the last five years (Bender, 2019).

The limitation of studies in the Portuguese language emphasizes the need to explore and enhance

the application of Open IE techniques in such language. However, works in this area publish their results as a package, making it difficult to use by non-computer experts. This article aims to contribute to this gap by presenting a hub of services that integrate methods from Open IE in Portuguese, enabling a single environment for users to extract triples and compare the results.

Our framework deserves a hub of Open IE services that implements each service as an Open IE method for the Portuguese language. We call this hub of Open IE services as FORMAS Open IE Framework.

This article is structured as follows: the next section describes the background. Section 3 presents each Open IE method. Section 4 describes our hub of Open IE services. Section 5 discusses our conclusions and some envisioning work.

2 Background

The advantages of Open IE, as outlined by (Gamallo, 2012), encompass domain independence, unsupervised extraction, and greater scalability. These features highlight Open IE’s adaptability, efficiency, and scalability in handling diverse subjects and extensive volumes of unstructured text.

Concerning architecture, the system’s architectural choice is pivotal in shaping Open IE systems’ development and maintenance. Various paradigms, including monolithic, microservices, layered, and event-driven architectures, offer distinct characteristics influencing scalability, flexibility, and modularity (Fowler, 2003; Newman, 2015a; Richardson, 2018).

A brief exploration of these architectural approaches follows:

Monolithic Architecture: Integrates all components into a single application, simplifying development and maintenance. While suitable for applications with precise requirements, it

080	may face scalability challenges as the system	DptOIE addresses specific linguistic construc-	129
081	grows (Fowler, 2003).	tions, such as:	130
082	<i>Microservices Architecture</i> : Divides a system	• Coordinative Conjunctions (CC): Handles	131
083	into independent services, promoting scala-	conjunctions like "and" or "or" to generate	132
084	bility and independent evolution. Facilitates	multiple triples from a single sentence.	133
085	modular development, updates, and mainte-	• Subordinate Clauses:Manages adjective, ad-	134
086	nance (Newman, 2015a).	verbial, and substantive clauses, linking them	135
087	<i>Layered Architecture</i> : Organizes the system	to the main clause to form coherent triples.	136
088	into abstraction levels, promoting modularity.	• Appositives: Derives additional triples from	137
089	While simplifying maintenance, inter-layer	sentences with appositives, creating synthetic	138
090	dependencies may limit scalability (Fowler,	clauses.	139
091	2003).	Consider the sentence "O diretor do hospital,	140
092	<i>Event-Driven Architecture</i> : Components com-	Júlio, vendeu sua fazenda." DptOIE extracts the	141
093	municate through asynchronous events, fa-	main triple: (O diretor do hospital; vendeu; sua	142
094	voring scalability and dynamic responsive-	fazenda). Moreover, it recognizes the appositive	143
095	ness. Allows distributed processing but re-	"Júlio" and generates a new triple: (O diretor do	144
096	quires careful management of asynchronous	hospital; é; Júlio). Additionally, it applies transi-	145
097	events (Richardson, 2018).	tivity to create an additional triple: (Júlio; vendeu;	146
098		sua fazenda).	147
099	The choice of architecture depends on project-	3.2 PTOiE-Flair	148
100	specific needs, with monolithic architectures of-	The PTOiE-Flair is a new OpenIE model based on	149
101	fering simplicity, microservices providing scalabil-	deep neural networks that enables the generation of	150
102	ity, layered architectures ensuring modularity and	triples given a sentence. PTOiE-Flair was trained	151
103	event-driven architectures supporting reactivity in	with two datasets, LSOI and S2, achieving SOTA	152
104	distributed systems. Each decision entails profound	results for the Portuguese language.	153
105	implications for system evolution and maintenance,	Consider the sentence "Os cachorros, que são	154
106	emphasizing the importance of considering project	mamíferos, são os melhores amigos do homem."	155
107	characteristics and requirements.	PortNOIE extracts two triples: ["Os cachor-	156
108	3 Services	ros"/ARG0, "são"/V, "mamíferos"/ARG1, "são"/V,	157
109	Services are implemented as Open IE methods. A	"os melhores amigos do homem"/ARG1], that is:	158
110	set of OpenIE methods for the Portuguese language	(i) <i>Os cachorros; são; os melhores amigos do</i>	159
111	was selected to be part of the first version of this	<i>homem</i> and (ii) <i>Os cachorros; são; mamíferos.</i>	160
112	hub of services: DptOIE, PTOiE-Flair, and Chat-	3.3 ChatGPT	161
113	GPT. We detailed each one as follows.	The use of ChatGPT for triple extraction repre-	162
114	3.1 DptOIE	sents an innovative approach to integrating natu-	163
115	DptOIE (Oliveira et al., 2023) is a method de-	ral language technologies with the extraction of	164
116	veloped for Open Information Extraction (OIE),	structured information. Considering the ChatGPT,	165
117	specifically designed for the Portuguese language.	carefully formulating prompts is essential to guide	166
118	The main objective of DptOIE is to extract valu-	the model in generating structured responses.	167
119	able information or "facts" from sentences by analyzing	Prompt example: "Provide a triple containing	168
120	their syntactic structure and dependencies. DptOIE	a subject, a relation, and an object based on the	169
121	has three main phases: Pre-processing, triple ex-	following statement: 'The event occurred when'."	170
122	traction and special cases.	We employed the ChatGPT as a service, pro-	171
123	The preprocessing carries out a <i>Tokenization</i> ,	viding custom prompts and receiving structured	172
124	<i>Part-of-Speech</i> (POS) tagging, and <i>Dependency</i>	responses as a triple structure. Initially, we utilized	173
125	<i>Analysis</i> to inputted sentences. DptOIE identifies	the davinci-003 variant of the ChatGPT model.	174
126	triples (subject, relation, object) by traversing the	However, in future deployments, new models such	175
127	dependency tree using a Depth-First Search (DFS)		
128	approach. Each triple consists of an argument		

as GPT-4 or GPT-4-turbo could be seamlessly integrated without compromising performance.

Consider the prompt: "Describe a situation in which" followed by a specific context. ChatGPT generates a structured response, such as "(a scientist; makes; a significant discovery)".

4 FORMAS OpenIE Framework

The monolithic architecture based on Django framework simplifies the development and maintenance of robust systems. According to (Newman, 2015b; Fowler, 2014), monolithic architectures provide an integrated approach, making implementing and managing functionalities easier.

Our architecture comprises two modules as depicted in Figure 1: a Frontend and a Backend. Our front end was integrated with Django, and our back end was initially developed as a service hub with three methods: Chatgpt, Ptoie-flair, and DptOIE.

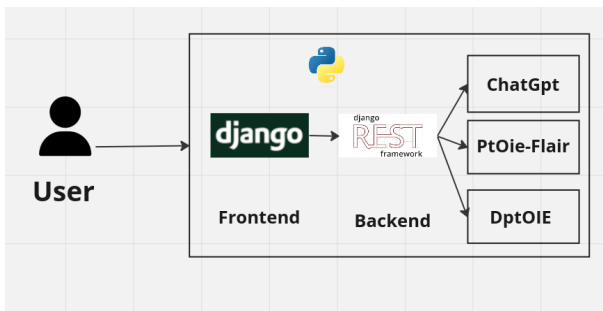


Figure 1: Architecture

Integration with Django enables the creation of a cohesive application, facilitating efficient communication between components. As mentioned by (Mokhtar, 2018), the use of monoliths is particularly advantageous when a pragmatic and efficiency-centered approach is sought, making it ideal for applications with clear and well-defined requirements.

In addition to the monolithic architecture, deploying the service on Cloud enhances resource adaptability for users. Utilizing cloud services provides dynamic scalability, allowing adjustments based on demand. According to (Kavis, 2014), Cloud services offers a reliable and flexible infrastructure, ensuring that users have access to resources tailored to their needs.

4.1 Graphical Interface

Our framework has a graphical interface depicted in Figure 2 that enables the use of visual elements to extract triples.



Figure 2: Interface

The user introduces a sentence in Portuguese, and our Framework answers with the triple extracted. The blue rectangle is omitted for anonymous review, and neither is the line under the title. Firstly, the user describes the sentence and selects which OpenIE method would prefer to extract the triples. Results are shown immediately after the button.

5 Conclusions and Future Work

The integration between the monolithic architecture using Django and deployment on the Cloud reinforces the robustness and adaptability of the Open Information Extraction services. This approach, supported by solid theoretical foundations, allows non-computer users to evaluate the extraction on three OpenIE methods.

In this work, while the primary goal was to provide a unified platform for users to access and utilize different OpenIE methods, the comparative analysis of these methods was not within the scope. Such comparisons entail diverse performance metrics and evaluation criteria, which could vary depending on the specific natural language processing tasks at hand. However, as future work, conducting comparative evaluations across various NLP tasks could offer valuable insights into the strengths and weaknesses of different extraction methods. By expanding the scope of evaluation beyond OpenIE,

242 we can further refine and optimize these methods to
243 better serve the needs of the Portuguese language
244 community. We envision enhancing the implementa-
245 tion and including more methods for the Open IE
246 Portuguese language community.

247 **Acknowledgments**

248 This material is partially based upon work sup-
249 ported by the FAPESB under grant INCITE
250 PIE0002/2022 and FAPESB TIC 0002/2015 and
251 CAPES Financial code 001.

252 **References**

- 253 Emily M. Bender. 2019. On the lack of study of non-
254 english languages in nlp. In *Proceedings of the First*
255 *ACL Workshop on Ethics in NLP*, pages 7–13.
- 256 Oren Etzioni, Michele Banko, Stephen Soderland, and
257 Daniel S. Weld. 2008. [Open information extraction](#)
258 [from the web](#). *Commun. ACM*, 51(12):68–74.
- 259 Martin Fowler. 2003. *Patterns of Enterprise Application*
260 *Architecture*. Addison-Wesley.
- 261 Martin Fowler. 2014. [Microservices](#).
- 262 Pablo Gamallo. 2012. Overview of open information
263 extraction. In *Open Information Extraction: Volume*
264 *377 of CEUR Workshop Proceedings*, pages 1–13.
- 265 Mike Kavis. 2014. *Architecting the Cloud: Design De-*
266 *isions for Cloud Computing Service Models*. John
267 Wiley & Sons.
- 268 Khaled Mokhtar. 2018. [The advantages and disadvan-](#)
269 [tages of monolithic and microservices architectures](#).
- 270 Sam Newman. 2015a. *Building Microservices: Design-*
271 *ing Fine-Grained Systems*. O’Reilly Media.
- 272 Sam Newman. 2015b. *Building Microservices: De-*
273 *signing Fine-Grained Systems*, 1st edition. O’Reilly
274 Media.
- 275 Leandro Oliveira, Daniela Barreiro Claro, and Marlo
276 Souza. 2023. [Dptoie: a portuguese open information](#)
277 [extraction based on dependency analysis](#). *Artif. Intell.*
278 *Rev.*, 56(7):7015–7046.
- 279 Chris Richardson. 2018. *Microservices Patterns: With*
280 *Examples in Java*. Manning Publications.