

Agricultural recommendation system for crop protection

Javier Lacasta, F. Javier Lopez-Pellicer, Borja Espejo-García, Javier Noguerras-Iso, F. Javier Zarazaga-Soria

Aragon Institute of Engineering Research (I3A), Universidad de Zaragoza, Spain

Abstract

Pests in crops produce important economic losses all around the world. To deal with them without damaging people or the environment, governments have established strict legislation and norms describing the products and procedures of use. However, since these norms frequently change to reflect scientific and technological advances, it is needed to perform a frequent review of affected norms in order to update pest related information systems. This is not an easy task because they are usually human-oriented, so intensive manual labour is required. To facilitate the use of this information, this work proposes the construction of a recommendation system that facilitates the identification of pests and the selection of suitable treatments. The core of this system is an ontology that models the interactions between crops, pests and treatments.

Keywords: Ontology creation, Ontology population, Data integration, Intelligent systems, Pest control

1. Introduction

Agriculture is a vital sector in the economy of any country, but depending on the crop between 26% and 80% of the agricultural production is lost because of pests (Oerke, 2006). Crop protection is vital but also challenging due to the multiple pests that affect them, such as insects, plant pathogens and weeds, and the toxic effects of most of the existing solutions (Alavanja, 2009). Because of

7 these effects, most countries have established strict regulations for their use and
8 promote non-chemical solutions (European Parliament, 2009).

9 In general, the norms about pest control are published in heterogeneous and
10 human oriented formats, so intensive manual labour is required to identify the
11 most suitable solution for a given pest. An example of this heterogeneity can be
12 found in the data collections provided by the Spanish Ministry of Agriculture¹
13 where the description of how to control each type of pest is distributed among
14 multiple heterogeneous textual sources. For example, each document has a lay-
15 out slightly different from the rest and the names of the pests in the document
16 title are variants of those used in the pest description. This lack of interoper-
17 ability affects critically tasks requiring some degree of data integration such as
18 identifying the different crops affected by a single organism, finding similitude
19 in the treatment of different species, and comparing the approved pesticides in
20 different countries. Additionally, as new products and techniques are frequently
21 approved, a continuous review is required (Ricci et al., 2010). This happens
22 not only in Spain, but also in many other countries such as United Kingdom²,
23 United States³ and Canada⁴.

24 To facilitate the usability of this information, we need systems able to provide
25 it in an integrated and harmonized way. For this task, in this paper, we propose
26 the “Pests in Crops and their Treatments” Ontology (PCT-O). To populate it,
27 we suggest a conversion process for the transformation of non-ontological het-
28 erogeneous resources into ontological ones. As use case, this process is applied to
29 transform content from selected Spanish data sources into instances according
30 to PCT-O model. Finally, we describe the structure of the information retrieval
31 (IR) system and the recommendation process that simplifies the identification
32 of a pest and the selection of a suitable treatment.

¹<http://www.mapama.gob.es/>

²<https://secure.pesticides.gov.uk/pestreg/>

³<https://www.epa.gov/pesticide-registration>

⁴<https://www.canada.ca/en/health-canada/services/consumer-product-safety.html>

33 2. State of the art

34 The use of ontologies is a classical solution to deal with heterogeneity and
35 interoperability problems. In the biology area, Walls et al. (2012a) remark how
36 semantic models facilitate the creation of intelligent applications that manage
37 living species information. The inference capability of ontologies are especially
38 relevant in the biology area, because it can be used in the taxonomic structures
39 used for classification to simplify conceptual interoperability, data integration
40 and search. However, the creation of ontologies is difficult. The main challenges
41 are the modelling of the information for the desired task, the availability of
42 data for population, and the data transformation complexity. Data modelling
43 is difficult due to different interpretations of the selected knowledge area. With
44 respect to data availability, the availability of data sources conditions the ex-
45 tension and depth of a semantic model. Something similar happens with data
46 transformation. Too complex or too heterogeneous data collections may not be
47 added to the model due to transformation costs.

48 Several works in the literature categorize living species, the interactions be-
49 tween them or the effects produced by chemical substances. This section de-
50 scribes the main works in these fields, remarks the parts of these models that
51 can be used to describe pest control information, and indicates the shortcomings
52 solved by the proposed PCT-O.

53 With respect to living being descriptions, the Integrated Information Tax-
54 onomic System (ITIS) (Integrated Taxonomic Information System, 2010) con-
55 tains taxonomic information of aquatic and terrestrial flora and fauna, the Cat-
56 alogue of Life model (Jones et al., 2000) describes 2 million of species, and the
57 NCBI taxonomy (Gene Ontology Consortium, 2004; Federhen, 2012) stores the
58 organism names and taxonomic lineages in the INSDC database. All these mod-
59 els provide a comprehensive collection of species but they do not provide very
60 detailed information about their features and behaviour. The search capabilities
61 of the portals providing them are limited to the use of names or database codes.

62 Other works provide extended taxonomies with additional information such

63 as species descriptions, biology, lifecycle, habitat, and interaction with other
64 species. An example of this type of works is Wikispecies (Wikimedia founda-
65 tion, 2017), which contains near half a million of species, although the informa-
66 tion provided for each species is limited. Focusing on plants, the U.S. plants
67 database (Natural resource conservation service, 2016) includes a quite detailed
68 textual description of U.S. plant, their distribution, life cycle, and common
69 pests. Another system is the European Nature Information System (EUNIS)
70 (Davies et al., 2004). It includes a large collection of species obtained from other
71 databases and indicates the geographical distribution and the level of extinction
72 threat of those species. A relevant work is the Encyclopedia of Life (Li et al.,
73 2004), which provides more detailed information about a million of species and
74 even a basic description of the interaction between species. However, it does
75 not detail the kind of interaction they have (predator, prey, symbiosis, and so
76 on). Sini (2009) describes the AGROVOC vocabulary, an agriculture thesaurus.
77 A part of it provides a taxonomy of living beings that includes the main used
78 crops and pests in the form of hierarchically related concepts. DBpedia (Auer
79 et al., 2007) also contains a formal structure for the information about living
80 species in Wikipedia and Wikispecies. However, the number of provided species
81 is more limited. Finally, GeoSpecies (DeVries, 2013) relates each concept to the
82 Encyclopedia of Life, Wikipedia, Wikispecies, NCBI, ITIS, and other similar
83 systems. Instead of providing proper information about the stored species, it
84 focuses on providing equivalences between the aligned models. The search capa-
85 bilities in these systems are more complete, allowing textual search in the data
86 content. In the semantic models, such as AGROVOC, DBpedia and GeoSpecies,
87 arbitrary searches are also possible.

88 Some works specifically focus on the interactions between species. Rodríguez-
89 Iglesias et al. (2017) propose an ontology that details the pathogens that affect
90 plants. It integrates data related to both plant physiology and plant pathology
91 with the objective of facilitating the interpretation of phenotypic responses and
92 disease processes. Similar to this, Walls et al. (2012b) analyse the infectious
93 diseases of plants and the pathogens that cause them. They reuse vocabular-

94 ies from other plant, pathogen and disease ontologies such as the Infectious
95 Disease Ontology (IDO) (Cowell and Smith, 2010). Finally, the Plant Ontology
96 Consortium (2002) defines a set of ontologies to describe plants, their genes, dis-
97 eases and growing process that include the relation between plants and harmful
98 virus and bacteria. All these models, as in the previous cases, provide semantic
99 searches that make possible detailed queries and precise results.

100 With respect to crop treatments, PubChem model (Fu et al., 2015) describes
101 chemical structures, biological activities and biomedical annotations. This in-
102 cludes pesticides and the environmental effects they produce. However, this in-
103 formation is text-based and it is not linked to any living species model. ChEBI
104 ontology is another model describing chemical substances (Degtyarenko et al.,
105 2008). It contains natural molecular entities and synthetic products that affect
106 living organisms. However, it also lacks a semantic relation with the species
107 affected by each chemical product. Here, depending on the part of the models,
108 textual or semantic searches are possible.

109 Other works integrate parts of all these and other agricultural aspects to-
110 gether. Damos (2013) proposes the definition of ontologies that allow describing
111 all the characteristics of cultivations. He also indicates the need to link the cre-
112 ated models to other related data collections that complement them. Damos
113 et al. (2017) show an ontology to describe pest and the treatments approved by
114 the Greek Ministry of Rural Development and Food. The core of the ontology
115 contains the pests that are related to the affected crops and existent treatments.
116 On a broader context, Athanasiadis et al. (2009) describe several ontologies for
117 data integration in the agricultural field. Especially relevant is their agricultural
118 activities ontology for crop management. Goumopoulos et al. (2009) describe
119 an ontology for precision agriculture. It focuses on describing plants and all the
120 technological and electronic devices that surround them in precision agricul-
121 ture. Finally, Rehman and Shaikh (2011) describe another precision agriculture
122 ontology whose core includes concepts for describing crops and their pests.

123 The objective of the ontology proposed in this paper (PCT-O) is to connect
124 crops, pests and treatments into a unified model. The formal description of liv-

125 ing species taxonomies can be managed with the previously described ontologies
126 such as NCBI taxon or GeoSpecies, the description of plant pathologies is cov-
127 ered by Rodríguez-Iglesias et al. (2017) illnesses ontology, and PubChem covers
128 the application of chemical substances. However, they do not model all the crop
129 protection aspects. Specifically, they do not cover the relation between crops,
130 pests that affect them, and the solutions approved by each country to deal with
131 them. Only Damos et al. (2017) make a proposal to relate information about
132 pests and treatments to the affected crops. However, they propose a high-level
133 model that does not provide detailed properties about each of the proposed
134 classes. The proposed PCT-O allows describing the conditions required by a
135 pest to produce outbreaks and the restrictions on the treatments.

136 **3. Structure of the PCT-O**

137 This section describes the ontology created for the description of pests, crops
138 and their treatments. The core of the proposed model can be considered as an
139 extension of the disease triangle described in Rodríguez-Iglesias et al. (2017),
140 which consists of a virulent pathogen, a susceptible host, and a propitious envi-
141 ronment. It has been extended to include non-pathogen pests and the definition
142 of treatments for the pests. We have also modelled the provenance of the in-
143 formation to allow updates and correction of errors in the sources and in the
144 generation process.

145 The ontology has been created with the Methontology methodology (Gómez-
146 Pérez et al., 2004). Specifically, the modelling has been guided to answer the
147 following competence questions: Which is the pest that is affecting a given crop?
148 Which treatment do I have to apply to deal with the pest? When do I have to
149 apply the treatment? What are the sanitary/environmental restrictions of the
150 treatment?

151 In the construction process of the PCT-O, we have put a special emphasis on
152 reusing existing models to improve the ontology interoperability. Specifically,
153 we have analysed widely used models of living species (which include both crops

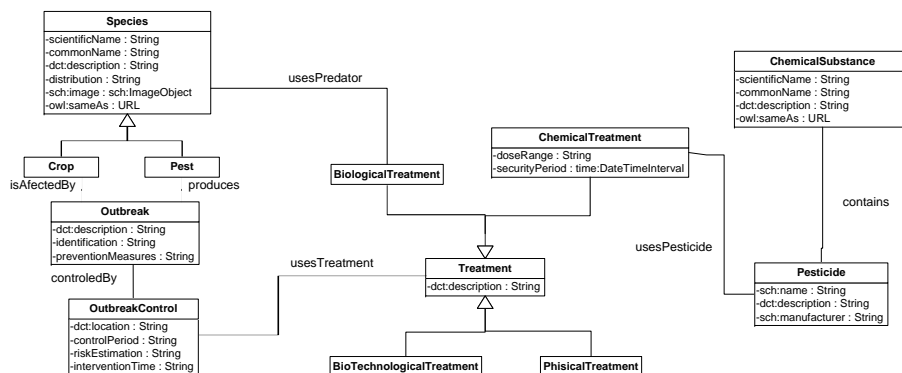


Figure 1: Plant affections and their treatment ontology

154 and pest) and chemical substances (which include pesticides) described in the
 155 state of the art section. The core *Species* and *ChemicalSubstance* classes in
 156 the model have DBpedia equivalents, and their instances are linked to NCBI
 157 taxon, PubChem, ChEBI ontology instances and the Spanish Wikipedia pages
 158 (using *owl:sameAs*). The connection between these elements has been guided
 159 according to the information provided in the Spanish guides for pest diagnosis
 160 and management.

161 The Spanish guides that detail the pest characteristics and treatments have
 162 provided us the terminology and relations used to construct the proposed on-
 163 tology. However, their lack of structure has forced us to use a coarse level of
 164 granularity for properties, leaving many of them as simple text fields. A finer
 165 granularity level is possible, but extracting the concepts and relations from
 166 the guides would require the definition of complex natural language processing
 167 (NLP) rules specific to each property. This issue is detailed in the discussion
 168 section.

169 Figure 1 shows the conceptual view of PCT-O. The main concept is the
 170 *Species* concept, which describes the name and characteristics of the included
 171 species. It has been specialized into *Crops* grown by farmers and *Pests* that
 172 harm the *Crops*. Crops that act as weeds can be classified as both types. The
 173 attributes are the common and scientific name the species, a description, its

174 distribution, images, and equivalency relations with other species models.

175 The *Outbreak* class models the interaction between crops and pests. It con-
176 tains a textual description of the produced symptoms, the identification and
177 analysis procedures used to establish that a pest is affecting a crop and the
178 existent prevention measures to reduce the risk of infection. It is based on the
179 IDO ontology, but our ontology also covers insects, plant pathogens and weeds.
180 It has been simplified because of the complexity of filling the description of
181 symptoms from the data sources.

182 The *OutbreakControl* class models the procedure to control a specific kind of
183 *Outbreak* and its location restrictions. Humidity and temperature are the main
184 triggers of outbreaks. Therefore, control procedures and recommendations may
185 vary depending on the climatology of each region. This class includes the period
186 of time in which the pest is harmful to the crop, the description of a way to
187 estimate the infection risk, the description of the best moment to take action
188 to reduce the damages, and the list of treatments approved in the location for
189 dealing with the pest.

190 The *Treatment* class describes four kinds of treatments: *Biological*, *Bio-*
191 *technological*, *Physical* and *Chemical*. Biological treatments make use of preda-
192 tors, physical treatments describe manual measures such as removing infected
193 fruits, bio-technological measures mostly use traps and pheromones, and chemi-
194 cal treatments use pesticides. Each treatment has a description of the treatment
195 itself. The chemical treatments are linked to the pesticides approved by the gov-
196 ernment (*Pesticide* class), the regulated amount and the legal period between
197 the application and the harvest.

198 The ontology describes the substances dangerous to the environment con-
199 tained in pesticides through the *ChemicalSubstance* class. It includes the com-
200 mon and scientific names of the substances and a description of the effects
201 caused and interactions with other species. We link the substances to Pub-
202 Chem, ChEBI ontology and the Spanish Wikipedia through the *owl:sameAs*
203 property. PubChem link is especially relevant as it contains information about
204 the environmental hazards produced by the chemical substances, and the rec-

205 ommended restrictions of use (e.g. many chemical substances must not be used
206 near water sources or some protected/commercial species). We think this infor-
207 mation is vital to be able to select appropriately the least aggressive solution
208 among the existent ones for a given place at a given time.

209 The ontology instances contain information extracted from multiple sources.
210 In this context, knowing the provenance of each piece of information is vital if
211 errors are detected or the sources change. Rodríguez-Iglesias et al. (2016) pro-
212 poses the use of a named graph structure in which the URI of the named graphs
213 are the base URI of the involved resources. We implement a similar solution by
214 using the PROV ontology (Lebo et al., 2013), which is recommended by W3C
215 for provenance description in the web. From PROV, we have used the *Bundle*
216 class and *hasDerivedFrom* property as our goal is to store the instance sources.
217 A *Bundle* is a named set of provenance descriptions that describe the common
218 provenance properties of a set of elements. *Bundles* contain the *hasDerivedFrom*
219 property that links the *Bundle* to the source file of the controlled elements. The
220 direct implementation of a *Bundle* is using a named graph. Named graphs define
221 collections of resources in a semantic repository under a single name and can be
222 annotated with the necessary properties. The combination of the *Bundles* pro-
223 vides the complete view of the provenance of the crops, pests and treatments.
224 Figure 2 shows an application example where the information extracted from
225 the “Agrotis Ipsillon” diagnosis guide is stored in a named graph and then inte-
226 grated with the rest of the instances for query. Since the information obtained
227 from each source is stored in different named graphs, it is possible to identify
228 their provenance by querying about the named graph that contains it.

229 **4. Ontology construction and population**

230 The backbone of the ontology instances are the NCBI taxon and the Spanish
231 Wikipedia for living species (crops and pests) and PubChem, ChEBI ontology,
232 and the Spanish Wikipedia for pesticide substances. The NCBI taxon, Pub-
233 Chem and ChEBI ontologies are well-known models in their respective fields

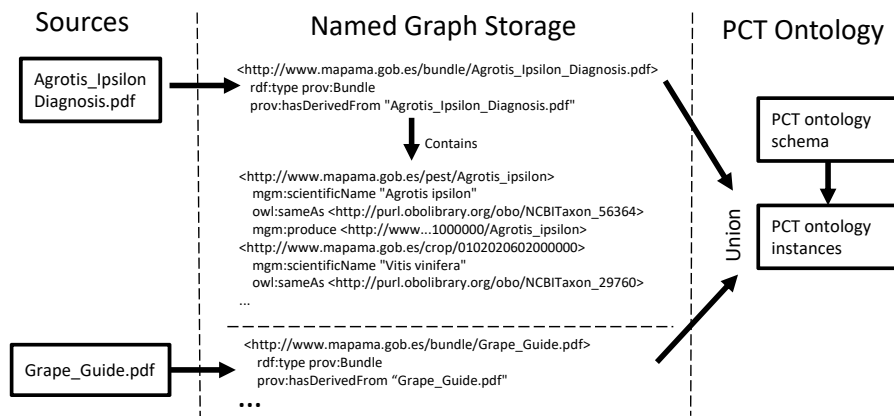


Figure 2: Example of provenance modelling

234 and provide the scientific names for each element (crop, pest and chemical sub-
 235 stances). Specifically, NCBI taxon provides a hierarchy of species useful for
 236 identification of families of crops. The Spanish Wikipedia provides alternative
 237 scientific and common names that are helpful in the disambiguation process.
 238 Each model has additional information about species and chemical substances
 239 such as taxonomic relations, definitions, chemical formula and so on. We do
 240 not currently use this information, but the linkage makes it accessible for future
 241 improvements.

242 To populate the PCT-O we have focused on the official information about
 243 crops and authorised pesticides maintained by the government of Spain. This
 244 section describes the data sources, the ontology construction and the process
 245 developed to extract the available information and represent it according to the
 246 ontology model.

247 4.1. Tools used for ontology construction

248 We have selected OWL (McGuinness et al., 2004) as the description model
 249 for our ontology and its instances. OWL is the most common RDF-based de-
 250 scription model in the semantic field and it enriches the description capabilities
 251 of RDF/RDFS (Brickley et al., 2014) by supporting complex relations between
 252 classes and detailed characterization of properties. The construction of the on-

253 tology has required the use of multiple tools and libraries to define the model
254 and populate it from the selected sources. The ontology has been created using
255 the Protégé editor⁵, a tool designed to facilitate the creation of OWL schemas.
256 With respect to the ontology population, it has required the extraction of infor-
257 mation from multiple PDF files. This has been done using Apache PDFBox⁶, a
258 Java library for PDF processing. For the processing of the extracted content, a
259 workflow that fills an Apache Jena⁷ triple-store (a RDF database that support
260 named graphs) has been created using Spring Batch⁸. Finally, the recommen-
261 dation tool is a very simple text interface that uses SPARQL (Prud et al., 2006)
262 (a language for querying RDF graphs) to extract the desired information from
263 the Jena triple-store.

264 *4.2. Data sources used for population*

265 The description of the effects that each pest has in each crop and the pro-
266 cesses established to detect and treat them have been obtained from the fol-
267 lowing heterogeneous document collections provided by the Spanish Ministry of
268 Agriculture: The laboratory diagnosis sheets of noxious species for crops created
269 by the phytosanitary diagnosis and survey laboratory, which is a collection of
270 464 scanned PDF documents describing plants, insects, bacteria and virus (sci-
271 entific and common names of the pests that affect crops, their distribution in
272 Spain, symptoms, detection measures and identification procedures); the guides
273 for the integrated control of pests created by the national plan for sustainable
274 use of pesticides, which is a collection of 21 digital PDF documents that describe
275 the crops affections in Spain and the recommendations for their treatment (com-
276 mon name of the crops, the common and scientific name of the noxious species,
277 control and prevention measures, and available non chemical treatments); and
278 the registry of pesticides approved by the national institute for agrarian research

⁵<https://protege.stanford.edu/>

⁶<https://pdfbox.apache.org/>

⁷<https://jena.apache.org/>

⁸<https://projects.spring.io/spring-batch/>

279 and technology, which is a repository containing 2375 PDF records detailing the
280 pesticides allowed in Spain, their composition and use restrictions.

281 The content of these sources connects the living species information with
282 the chemical substances used on them. The main issue of these collections
283 is their heterogeneity. None of these data sources is completely structured and
284 uniform. Some parts have a tabular structure, but most of them are described as
285 paragraphs of plain text. The text sections are similar between documents but
286 not exactly equivalent. Additionally, the quality of several scanned documents
287 is low, making data extraction difficult.

288 *4.3. Population process*

289 We have followed the population process described in Figure 3. The first step
290 has been to extract the textual content and available images from the source
291 PDF files. Then, each type of source has been parsed to identify the elements
292 required in the ontology. Textual content is used for filling the different proper-
293 ties of the instances, while the images are stored as a graphical representation
294 of each concept. All the extracted images are stored, independently of the rele-
295 vance of their content. To simplify data integration, each extracted resource is
296 aligned to the previously described ontologies using the common and scientific
297 name of crops, pests and chemical substances as matching text. Having identi-
298 fied the species/chemical substances in the resources, their integration is direct.
299 The first half of the process is dependent of the selected sources, but the second
300 half can be directly used for integrating future additional data collections.

301 In the data extraction step, if the origin of the PDF file is analogical (scan-
302 ning of a printed document), the OCR process in the PDFBox library is applied
303 to extract the text. However, scan quality of the source files limits the quality
304 of the extracted content. Most of the extracted text contains minor errors due
305 to bad recognition of some characters, but a few have higher error rates. In
306 addition to this, the non-plain text parts of the documents are not correctly
307 extracted due to PDFBox limitations (e.g., captions of photos or tabular infor-
308 mation).

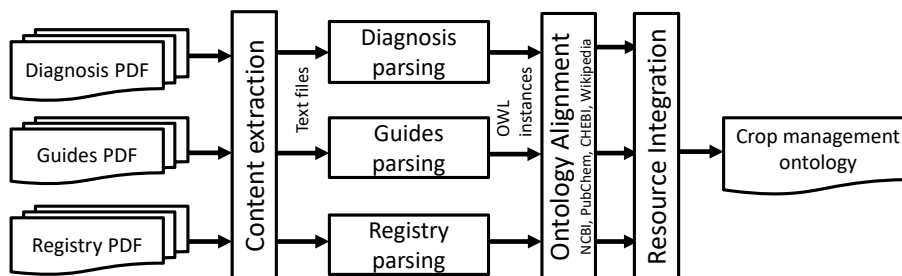


Figure 3: Ontology population process

309 The parsing step makes use of the fact that all the analysed sources are
 310 divided into sections whose content mainly corresponds with properties of the
 311 defined model. It identifies these sections according to a list of predefined head-
 312 ers for each type of document that contain all the variant forms found for the
 313 sections names and structure of the source documents. Additionally, we have
 314 defined specific rules containing syntactic patterns describing textual construc-
 315 tions in the documents when describing the common or scientific name of a
 316 species. The extracted information and its provenance information is stored
 317 according to the PCT-O model.

318 The alignment step matches the extracted resources describing species (crops
 319 and pests) with the NCBI taxon and the Spanish Wikipedia, and the chemical
 320 substances with respect to the Spanish Wikipedia, PubChem, and ChEBI on-
 321 tologies. The alignment of the species is used to directly merge the information
 322 of the involved data collections. The alignment of the chemical substances is
 323 used to facilitate the identification of equivalences between the different prod-
 324 ucts used to deal with the pests.

325 The alignment has been performed looking for equivalences in the scientific
 326 names of species and chemical substances contained in the documents. The
 327 complexity of this alignment process has come from the need of identifying and
 328 correcting the errors in the sources, and because of the existence of synonyms
 329 and variants of names of the living beings and chemical substances. To deal with
 330 these problems, we have performed the following alignment sub-steps. First, we

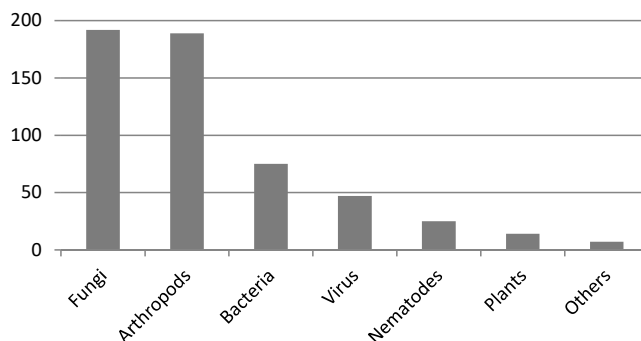


Figure 4: Classification of pests

331 have extended the available synonyms and variant names for each extracted
 332 crop/pest with additional names obtained from the Spanish Wikipedia. This
 333 has been done looking for the common names in the Spanish Wikipedia and
 334 extracting the scientific ones contained in the corresponding info-boxes. Then,
 335 all the scientific names are matched (exact match) with the corresponding on-
 336 tology/model (NCBI, PubChem, ChEBI). If a match is found, the alignment is
 337 established. If there is no correspondence, we have used the Levenshtein dis-
 338 tance (Levenshtein, 1966) to identify matches with minor errors and variants of
 339 the scientific names. For this comparison, the scientific names are normalized
 340 removing abbreviations, numbers, and texts in brackets. Name heterogeneity
 341 has led us to use a threshold of 20% of the name size to decide if the most similar
 342 name can be aligned or not. Therefore, shorter names allow smaller differences
 343 than longer ones. This threshold has been selected experimentally to reduce the
 344 number of incorrectly aligned concepts (we prefer to leave them unaligned).

345 The resulting ontology consists of 549 pests that affect 462 crops through
 346 3471 outbreaks. Figure 4 shows the pests in the model aggregated by family.
 347 It can be observed that most of them are fungi and arthropods. In addition to
 348 those, there are virus, bacteria, nematodes and other plants. A few pests are
 349 from species that do not fit in the previous categories. To deal with these pests,
 350 there are 42397 different chemical treatments involving 2109 pesticides with 566
 351 different chemical substances, and 219 alternative treatments.

352 A manual review of the ontology has shown that 96.12% of the species (pests
353 and crops) have been correctly aligned to their scientific name in NCBI Ontol-
354 ogy. The main source of errors are problems in the description of the names of
355 the sources (e.g., “summer cereals”), the use in the sources of the fruit name
356 instead of the plant name or the lack of equivalences for some of the used com-
357 mon names. We have also reviewed the quality of the extracted description of
358 the species, the symptoms and the information related to prevention and inter-
359 vention time. Here the quality is worse due to the difficulty of extracting the
360 content. There are almost no records without syntactic errors. Most of them are
361 small, but to be usable, it is required to correct them through a manual proof-
362 reading. Something similar happens with treatments: the extracted information
363 has been correctly assigned to the corresponding concepts in the ontology, but
364 there are many syntactic errors caused by the extraction. Finally, we have also
365 reviewed the alignment of the chemical substances with the ChEBI database
366 (PubChem is linked to it). The result shows that just 59.9% of the chemical
367 substances have been correctly aligned, 27.7% of them are left unaligned and
368 the rest (12.4%) are incorrectly aligned. This alignment problem is caused by
369 the lack of correspondence between the Spanish common/scientific names for
370 the chemical substances in the sources and the Spanish Wikipedia. The Span-
371 ish Wikipedia has proven to be a good source to align common and scientific
372 names of living species but its coverage for chemical substances is much worse.
373 It does not describe many specific substances, thus the Spanish names cannot
374 be aligned to the English ones in the selected ontologies.

375 From these data, it can be observed that current crop protection is com-
376 pletely focused around the use of chemical products. There are many more
377 chemical solutions than alternative ones, and their amplitude of action is also
378 broader because they affect several pests. With respect to alternative ap-
379 proaches, they are only able to deal with a small set of the pests (mainly insects)
380 but they do not have secondary effects for humans or nature.

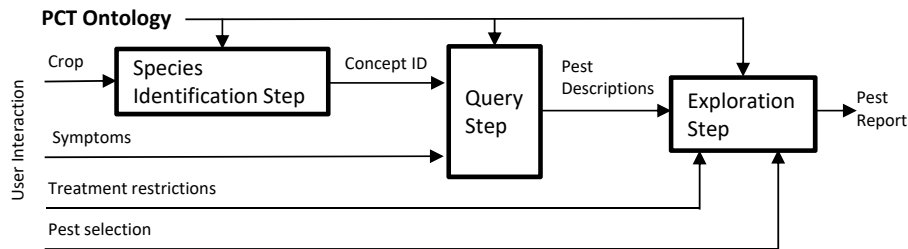


Figure 5: Query process

381 4.4. Recommendation system scenario

382 This section describes the developed IR-based recommendation system, con-
 383 structed on top of PCT-O to obtain directly complex information useful for crop
 384 protection, and describes its potential and limitations. Figure 5 shows the dif-
 385 ferent components of this process. These components use SPARQL queries to
 386 process the ontology and provide the results. The species identification step
 387 finds the crop concepts that correspond with the ones used in the query. Here,
 388 all the registered variants of common and scientific names are matched with the
 389 query term and the concept that matches it is returned. The query step identi-
 390 fies the pests that affect a crop with the symptoms indicated by the user. Since
 391 the species are defined in a taxonomical way and several of the relations are at
 392 category level (e.g., citric or fungus), any search by a member of these categories
 393 can be expanded to obtain all the pests affecting to its category. Finally, the ex-
 394 ploration step starts when the user selects a pest from the set obtained through
 395 the query step. Then, the local pest information and treatments are selected
 396 based on the user restrictions. If information from additional countries were
 397 added, it would be also possible to restrict solutions for products cultivated for
 398 exportation or even to identify better solutions than the one currently approved
 399 in the residing country.

400 Because of the coarse granularity level of the ontology, the query and explo-
 401 ration restrictions have to be done on text fields. This is a system limitation as
 402 text match solutions have problems related to synonymy, polysemy and multi-
 403 ple variant forms that reduce match quality. In this system, we have not used

404 provenance information, because their main purpose is for tasks related to model
405 updates, and versioning.

406 This recommendation system shows how PCT-O facilitates identification
407 tasks, but PCT-O also allows direct queries to list all the available treatments
408 for a pest in a crop. In this case, there is no ambiguity problems because it is
409 a direct query about specific elements that are perfectly identified.

```
1. Query = Crop:"Lemon tree", Symptoms:"Brown leaves", Treatments:"Biological"
2. Species identification step:
Select ?crop where {{?crop mgm:scientificName ?name. FILTER regex(?name, "Lemon tree", "i" )}
    union {?crop mgm:commonName ?name. FILTER regex(?name, "Lemon tree", "i" )}
    Result: http://www.mapama.gob.es/crop/0102020104000000 <- Citrus Limon URI
3. Query step:
Select ?outb where {{<http://www.mapama.gob.es/crop/0102020104000000> mgm:isAffectedBy ?outb}
    union {<http://www.../0102020104000000> skos:broader+ ?crop. ?crop mgm:isAffectedBy ?outb}.
    ?outb dc:description ?descr. FILTER regex(?descr, "Brown leaves", "i" )}
    Result: http://www.mapama.gob.es/ourbreak/0102020100000000/Tetranychus_urticae
        http://www.mapama.gob.es/ourbreak/0102020100000000/Citrus_exocortis_viroid_(CEVd)
4. Exploration step:
Select ?treatment where {<http://www.../Tetranychus_urticae> mgm:isControlledBy ?control.
    ?control mgm:usesTreatment ?treatment.
    ?treatment rdf:type <http://www.mapama.gob.es/vocabulary#BiologicalTreatment>}
```

Figure 6: Example of query specification and SPARQL queries performed

410 As a summarised example of this IR flow, we describe how the query de-
411 picted in Figure 6 is executed (it is simplified and just the concept identifier is
412 returned). The current query interface allows introducing the query terms to
413 search in the crop name, symptoms produced by the pest, and restrictions in
414 the treatment. The selected query (1) searches for a pest affecting the “Lemon
415 tree” that produces “Brown leaves” and how to treat it with a biological treat-
416 ment. The species identification step (2) directly matches the “Lemon tree”
417 species name with the “Citrus limon” concept in the ontology. “Citrus limon”
418 has no direct specification of pests as they are common to all ”Citrus” family.
419 Thus, the query step (3) expands the query to the ”Citrus” species and finds
420 two different pests, “Citrus exocortis viroid” and “Tetranychus urticae”, that
421 produce “Brown leaves”. For this expansion, we use a crop taxonomy extracted

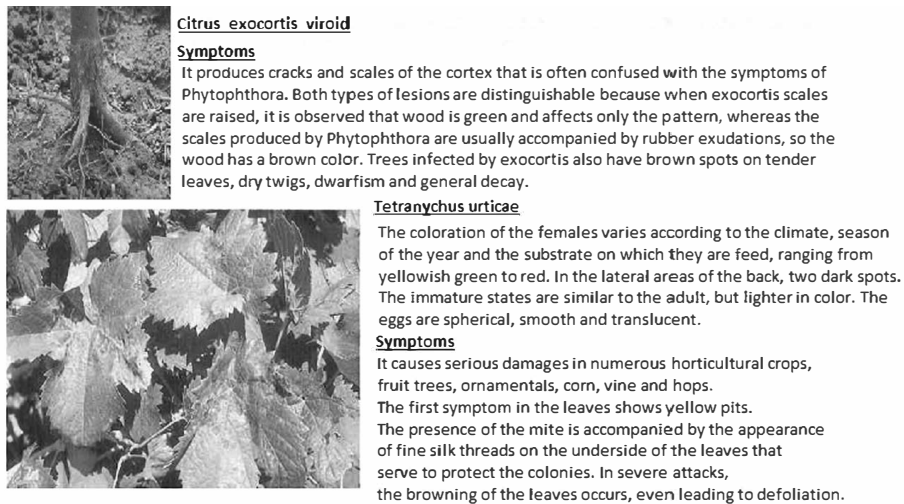


Figure 7: Example of information returned by the Query Step

422 from the sources, but since NCBI is linked to the concepts, it also could be used
 423 for this task. Figure 7 shows a composition of the information that can be
 424 returned in the Query Step (the original Spanish text has been translated to
 425 English to facilitate its understanding). Finally, given the “Tetranychus ur-
 426 ticae”, the exploration step (4) returns the available biological treatments for
 427 it, which consists in releasing predators such as Amblyseius (Neoseiulus) cali-
 428 fornicus, Phytoseiulus persimilis and Diptera Feltiella acarisuga.

429 Two problems have been found in this query system. First, source infor-
 430 mation is sometimes imprecise or incomplete. This is the case of the “Citrus
 431 exocortis viroid” that has no description. This lack of information can limit the
 432 ontology usability. The second issue is related to the generality of the infor-
 433 mation. For species that attack multiple crops, sources only provide the most
 434 general and representative examples. In this case, the “Citrus exocortis viroid”
 435 image is focused on roots, because the main symptom focuses there (leaves col-
 436 oration is secondary). In the “Tetranychus urticae” case, the image shows a
 437 leaf affected by the pest, but from a plant different from the “Citrus limon”.
 438 Correcting both issues would require to increase the amount and precision of

439 the data sources available.

440 5. Discussion

441 As indicated in the state of the art section, there are several models for
442 the description of species and chemical substances, but only Damos (2013) and
443 Damos et al. (2017) provide some relation between crops, pests, and treatments.
444 PCT-O goes a step further by including the description of the conditions of these
445 relations. Therefore, in PCT-O, it is possible to specify the period of time when
446 a pest is harmful, when it is needed to react, and the nature of the treatments.
447 PCT-O also includes provenance information to keep track of the data sources.
448 The next closest solution is the PubChem database (and ontology) that de-
449 scribes thousands of chemical substances and their application in the industry.
450 For the appropriate substances, it indicates the common name of the crops to
451 which the substance can be applied according to USA legislation. However, it
452 is not linked to any species ontology and may be ambiguous. Additionally, it
453 indicates neither a detailed list of the noxious species the chemical substance
454 can deal with, nor the symptoms, periods of control or chemical alternatives.

455 In the analysed scenario, we have shown how PCT-O helps in terms of
456 interoperability and data integration between crops, pests and treatments in-
457 formation. Thanks to it, it is possible to construct a semantic recommendation
458 system that helps to determine the pests that affect each crop and how to treat
459 them. The crops, pests, and pesticides are linked to commonly used ontologies
460 and taxonomies. This removes name ambiguity and allows comparing solutions
461 adopted in different regions or countries.

462 The population of the ontology with Spanish official data has illustrated the
463 complexity of obtaining a complete model from the available official sources.
464 Data quality has been an issue that has complicated the data transformation
465 and it has added errors. We have found several cases where a correct equiva-
466 lence has not been found and chemical substances have been incorrectly aligned.
467 The cause of this is mainly due to the incompleteness of Spanish Wikipedia

468 in biology/chemistry area and the similarity between some scientific names of
469 species/chemical substances. Another identified issue is related to the com-
470 pleteness and overlap of the data sources. Each data source was created by its
471 producer with a different purpose and they do not completely overlap. For in-
472 stance, the guides only cover a subset of species described in the diagnosis files.
473 As a result, the populated ontology does not have a uniform coverage: some
474 species are very detailed, other ones contain very limited information. These
475 restrictions reduce the usability of the extracted information, but it is a good
476 starting point for future improvement.

477 Because of the automatic nature of the population process and the hetero-
478 geneity of the sources, the resulting collection requires manual validation. For
479 this task, the stored provenance information becomes vital as incorrect or poorly
480 described instances can be traced to the original sources, allowing the detection
481 of the source documents with errors, so they can be fixed.

482 Although we have focused on Spain data for the population step, information
483 from other countries could be added. Countries such as U.S., United Kingdom
484 or Canada also provide the information required to populate this ontology in
485 heterogeneous formats, but specific extraction and transformation steps for each
486 new source format would be required. The step that align each species/chemical
487 with the selected ontologies and the final integration phase could be reused.

488 A limitation of PCT-O is the selected semantic granularity of the model. The
489 information contained in fields such as pest description, control period, identi-
490 fication procedures, or intervention time is described as plain text, so queries
491 on these fields are imprecise. For example, when querying for “Brown leaves”
492 as pest symptom, pests that only produce brown leaves in some specific situa-
493 tions will be returned with the same importance than pests with brown leaves
494 as representative symptom. Solving this problem would require to extend the
495 ontology to allow a precise description of such content. However, available in-
496 formation is so heterogeneous that cannot be automatically interpreted only
497 with the information contained in the source files. For example, in the period of
498 control of a crop, it is important to consider the growth stage, temperature and

499 humidity. The growth state is sometimes properly described (e.g., flowering),
500 but other times it is referenced using periods of months or seasons (e.g., May).
501 This must be interpreted depending on the place and the climate conditions
502 of a given year. The same happens with the humidity or temperature. Some
503 descriptions are quite clear (e.g., temperature under 25 degrees), but others
504 need human interpretation (e.g., high temperature). In this context, a semantic
505 baseline for each crop must be defined to allow the mapping of all the imprecise
506 descriptions to measurable values. We have done a preliminary processing to
507 identify the common temperature and humidity patterns in the source docu-
508 ments and more than 80 different rules have been needed. Additionally, we had
509 to perform approximations that are crop and pest dependent. For instance,
510 many documents say that a crop is vulnerable to a pest with high temperature,
511 but how much temperature is “high”? To model it semantically, this must be
512 translated to a numerical range (as it is in many other descriptions). However,
513 with the source information alone it is not possible to determine a precise value,
514 and an approximation must be given. Due to these approximations, we think
515 that the fine grain semantic extraction can only be useful as an initial step in
516 IR process. The final decision must be taken by the user who has interpret the
517 original description.

518 **6. Conclusions**

519 This work proposes the PCT-O ontology, a model to describe the outbreaks
520 that pests produce to crops and the approved ways to treat them. Currently,
521 there are several ontologies to describe taxonomies of living beings but none
522 allows describing their inter-relations as the PCT-O ontology. As use case for
523 this ontology, we propose a recommendation system that helps to identify the
524 pests affecting a crop and their treatments.

525 The ontology has been populated with official information in Spain about
526 crops, pests and approved treatments. This process has been complex due to
527 the heterogeneity, format and quality of the data sources. The extraction and

528 source errors, complemented with synonymy and name variants, have forced us
529 to use a disambiguation process of scientific names based on the alignment of
530 species and chemical substance records with ontologies such as NCBI, PubChem,
531 ChEBI and Wikipedia. The resulting model has been tested in a suggestion use
532 case to determine how to identify a pest and select a treatment. Additionally,
533 it can be used for tasks such as the identification of outbreaks, identification of
534 location-based related conflicts with the treatments, and comparison of solutions
535 between country legislations.

536 A first area of future work is to integrate treatments adopted by other coun-
537 tries for the same illnesses/pests in the population of the ontology. This will
538 require extending the extraction and parsing step to deal with the additional
539 data sources, but it will allow complementing the pest descriptions and compar-
540 ing the approved treatments to detect differences between regions. These
541 differences may show gaps in country legislations, and allow identifying better
542 solutions for a region than the currently approved ones.

543 Another interesting extension would be to include other aspects of the use
544 of chemical substances in the land. For example, PubChem repository contains
545 information about the hazards of the use of the chemical substances, such as
546 “Very toxic to aquatic life with long lasting effects”). This information merged
547 with water flow, crops or protected species distribution maps can be useful to
548 determine the areas where a product can be used, or suitable alternatives for
549 areas that forbid it. A complementary source of this information is the EU - Pes-
550 ticide Database (European Commission, 2005) that stores the list of substances
551 approved in each European member state for their use as pesticides. Finally, the
552 ontology could be extended to integrate more detailed information about crops
553 and their varieties. For example, the Spanish Ministry of Agriculture provides a
554 collection of descriptive sheets containing information about the different crop
555 varieties used in Spain. This collection provides information about the growth
556 conditions, performance and resistance of the different varieties of species. This
557 could be used to recommend the best variety for a field given its climate and
558 the distribution of the registered pests.

559 **Acknowledgments**

560 This work has been partially supported by the Spanish Government through
561 the project TIN2017-88002-R. The work of Borja Espejo-Garcia has been par-
562 tially supported by Aragón Government through the grant number C38/2015.

563 **References**

- 564 M. C. R. Alavanja. Pesticides use and exposure extensive worldwide. *Reviews*
565 *on Environmental Health*, 24(4):303–309, 2009.
- 566 I. N. Athanasiadis, A. E. Rizzoli, S. Janssen, E. Andersen, and F. Villa. On-
567 tology for seamless integration of agricultural data and models. In *Conf. on*
568 *Metadata and Semantic Research*, pages 282–293, 2009.
- 569 S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives. Dbpedia:
570 A nucleus for a web of open data. *The semantic web*, pages 722–735, 2007.
- 571 D. Brickley, R. V. Guha, and B. McBride. RDF Schema 1.1. *W3C recommen-*
572 *dation*, 2014.
- 573 L. G. Cowell and B. Smith. Infectious disease ontology. In *Infectious disease*
574 *informatics*, pages 373–395. Springer, 2010.
- 575 P. Damos. Semantics and emergent web-3 technologies: Modern challenges
576 for integrated fruit production systems towards internationalization. *IOBC-*
577 *WPRS Bulletin*, 91:133–142, 2013.
- 578 P. Damos, S. Karampatakis, and C. Bratsas. Representing and integrating agro
579 plant-protection data into semantic web through a crop-pest ontology: The
580 case of the greek ministry of rural development and food (GMRDF) ontology.
581 *IOBC-WPRS Bulletin*, 123:122–127, 2017.
- 582 C.E. Davies, D. Moss, and M. O. Hill. EUNIS habitat classification. Technical
583 report, European Environment Agency-European Topic Centre on Nature
584 Protection and Biodiversity, 2004.

- 585 K. Degtyarenko, P. de Matos, M. Ennis, et al. ChEBI: a database and ontology
586 for chemical entities of biological interest. *Nucleic acids research*, 36(1):344–
587 350, 2008.
- 588 P. J. DeVries. GeoSpecies knowledge base, 2013.
- 589 European Commission. EU pesticides database. Online database, 2005.
- 590 European Parliament. Regulation (EC) 1107/2009 of the European Parliament
591 and of the Council. Technical report, EU, 2009.
- 592 S. Federhen. The NCBI taxonomy database. *Nucleic acids research*, 40(1):
593 136–143, 2012.
- 594 G. Fu, C. Batchelor, M. Dumontier, J. Hastings, E. Willighagen, and E. Bolton.
595 PubChemRDF: towards the semantic annotation of PubChem compound and
596 substance databases. *Journal of Cheminformatics*, 7(34), 2015.
- 597 Gene Ontology Consortium. The gene ontology (go) database and informatics
598 resource. *Nucleic acids research*, 32(1):258–261, 2004.
- 599 A. Gómez-Pérez, M. Fernández-López, and O. Corcho. *Ontological Engineering*,
600 chapter Methodologies and methods for building ontologies. Methontology,
601 pages 125–142. Advanced Information and Knowledge Processing. 2004.
- 602 C. Goumopoulos, A. D. Kameas, and A. Cassells. An ontology-driven system
603 architecture for precision agriculture applications. *International Journal of*
604 *Metadata, Semantics and Ontologies*, 4(1-2):72–84, 2009.
- 605 Integrated Taxonomic Information System. Integrated taxonomic information
606 system on-line database, 2010.
- 607 A. Jones, X. Xu, N. Pittas, et al. Spice: A flexible architecture for integrating
608 autonomous databases to comprise a distributed catalogue of life. In *Int.*
609 *Conf. on Database and Expert Systems Applications*, pages 981–992, 2000.

- 610 T. Lebo, S. Saho, and D. McGuinness. PROV-O: The PROV ontology. Recommendation, W3C, April 2013.
- 611
- 612 V. I. Levenshtein. *Binary codes capable of correcting deletions, insertions, and*
613 *reversals*, volume 10, pages 707–710. 1966.
- 614 W. Li, R. W. Byrnes, J. Hayesa, et al. The encyclopedia of life project: grid
615 software and deployment. *New Generation Computing*, 22(2):127–136, 2004.
- 616 D. L. McGuinness, F. Van Harmelen, et al. OWL web ontology language
617 overview. *W3C recommendation*, 2004.
- 618 Natural resource conservation service. The plants database, 2016.
- 619 E. C. Oerke. Crop losses to pests. *Journal of Agricultural Science*, 144(31-43),
620 2006.
- 621 Plant Ontology Consortium. The plant ontology consortium and plant ontolo-
622 gies. 3(2):137–142, 2002.
- 623 E. Prud, A. Seaborne, et al. SPARQL query language for RDF. 2006.
- 624 A. Rehman and Z. Shaikh. Ontagri: scalable service oriented agriculture on-
625 tology for precision farming. In *Int. Conf. on agricultural and biosystems*
626 *engineering vols*, pages 1–2, 2011.
- 627 P. Ricci, M. Barzman, F. Bigler, et al. Integrated pest management in europe.
628 Technical report, ENDURE Network, 2010.
- 629 A. Rodríguez-Iglesias, A. Rodríguez-González, A. Irvine, et al. Publishing fair
630 data: an exemplar methodology utilizing phi-base. *Frontiers in plant science*,
631 7:641, 2016.
- 632 A. Rodríguez-Iglesias, M. Egana Aranguren, A. Rodríguez-González, and M. D.
633 Wilkinson. Plant-pathogen interactions ontology (PPIO). In *Int. Conf. on*
634 *Bioinformatics and Biomedical Engineering*, 2017.

- 635 M. Sini. Semantic technologies at FAO. *Agricultural information management*
636 *standards, International Society for Knowledge Organization (ISKO)*, 3, 2009.
- 637 R. L. Walls, B. Athreya, L. Cooper, et al. Ontologies as integrative tools for
638 plant science. *American journal of botany*, 99(8):1263–1275, 2012a.
- 639 R. L. Walls, B. Smith, J. Elser, et al. A plant disease extension of the infectious
640 disease ontology. In *ICBO*, pages 1–5, 2012b.
- 641 Wikimedia foundation. Wikispecies: free species dictionary, 2017.