



HAL
open science

PGxO and PGxLOD: a reconciliation of pharmacogenomic knowledge of various provenances, enabling further comparison

Pierre Monnin, Joël Legrand, Graziella Husson, Patrice Ringot, Andon Tchechmedjiev, Clement Jonquet, Amedeo Napoli, Adrien Coulet

► To cite this version:

Pierre Monnin, Joël Legrand, Graziella Husson, Patrice Ringot, Andon Tchechmedjiev, et al.. PGxO and PGxLOD: a reconciliation of pharmacogenomic knowledge of various provenances, enabling further comparison. BMC Bioinformatics, 2019, 20 (S4), 10.1186/s12859-019-2693-9 . hal-02103899

HAL Id: hal-02103899

<https://inria.hal.science/hal-02103899v1>

Submitted on 19 Apr 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

METHODOLOGY

Open Access



PGxO and PGxLOD: a reconciliation of pharmacogenomic knowledge of various provenances, enabling further comparison

Pierre Monnin^{1*}, Joël Legrand¹, Graziella Husson¹, Patrice Ringot¹, Andon Tchechmedjiev², Clément Jonquet^{2,3}, Amedeo Napoli¹ and Adrien Coulet^{1,3}

From The 2017 Network Tools and Applications in Biology (NETTAB) Workshop
Palermo, Italy. 16–18 October 2017

Abstract

Background: Pharmacogenomics (PGx) studies how genomic variations impact variations in drug response phenotypes. Knowledge in pharmacogenomics is typically composed of units that have the form of ternary relationships *gene variant – drug – adverse event*. Such a relationship states that an adverse event may occur for patients having the specified gene variant and being exposed to the specified drug. State-of-the-art knowledge in PGx is mainly available in reference databases such as PharmGKB and reported in scientific biomedical literature. But, PGx knowledge can also be discovered from clinical data, such as Electronic Health Records (EHRs), and in this case, may either correspond to new knowledge or confirm state-of-the-art knowledge that lacks “clinical counterpart” or validation. For this reason, there is a need for automatic comparison of knowledge units from distinct sources.

Results: In this article, we propose an approach, based on Semantic Web technologies, to represent and compare PGx knowledge units. To this end, we developed PGxO, a simple ontology that represents PGx knowledge units and their components. Combined with PROV-O, an ontology developed by the W3C to represent provenance information, PGxO enables encoding and associating provenance information to PGx relationships. Additionally, we introduce a set of rules to *reconcile* PGx knowledge, i.e. to identify when two relationships, potentially expressed using different vocabularies and levels of granularity, refer to the same, or to different knowledge units. We evaluated our ontology and rules by populating PGxO with knowledge units extracted from PharmGKB (2701), the literature (65,720) and from discoveries reported in EHR analysis studies (only 10, manually extracted); and by testing their similarity. We called PGxLOD (*PGx Linked Open Data*) the resulting knowledge base that represents and reconciles knowledge units of those various origins.

Conclusions: The proposed ontology and reconciliation rules constitute a first step toward a more complete framework for knowledge comparison in PGx. In this direction, the experimental instantiation of PGxO, named PGxLOD, illustrates the ability and difficulties of reconciling various existing knowledge sources.

Keywords: Knowledge engineering, Knowledge comparison, Semantic web, Ontology, Pharmacogenomics, Linked open data

*Correspondence: pierre.monnin@loria.fr

¹Université de Lorraine, CNRS, Inria, LORIA, 54000 Nancy, France
Full list of author information is available at the end of the article



Background

In this article, we present a simple ontology named PGxO, developed to reconcile and trace knowledge in pharmacogenomics (PGx). We instantiated this ontology with knowledge of various origins to both illustrate the relevance of the ontology and constitute a Linked Open Data (LOD) data set for PGx [1].

PGx itself studies how genomics impact individual variations in drug response phenotypes [2]. Knowledge in pharmacogenomics is of particular interest for the implementation of personalized medicine, i.e. a medicine tailoring treatments (chosen drugs and dosages) to every patient, in order to reduce the risk of adverse events. Indeed, best known examples of PGx knowledge already led to the development of clinical guidelines and practices [3] that recommend considering individual genotype when prescribing some particular drugs such as abacavir (an anti-retroviral) or fluorouracil (an anti-neoplastic) [4, 5].

Units of PGx knowledge typically have the form of a ternary relationship *gene variant – drug – adverse event* stating that a patient having the gene variant and being treated with the drug will have a higher risk of developing the mentioned adverse event. For example, one relationship is *G6PD:202A – chloroquine – anemia*, stating that patients with the 202A version of the G6PD gene and treated with chloroquine (an anti-malarial drug) are likely to experience anemia (an abnormally low level of red blood cells in blood).

An increasing volume of state-of-the-art knowledge in PGx can be found in reference databases, such as PharmGKB [6], or in the biomedical literature [7]. But, a large part of this knowledge is suffering from a lack of validation, or “clinical counterpart” [8], and is not yet ready to be translated into clinical guidelines and practices. For example, about 91% (as of July 2018) of the relationships listed in PharmGKB are qualified with a level of evidence 3 or 4, corresponding, at best, to an unreplicated study or to multiple studies that show a lack of evidence for the relationship [6]. On the other hand, PGx knowledge can also be discovered from clinical data, such as Electronic Health Records (EHRs), particularly when those are linked to DNA biobanks [9–11]. In this case, discovered knowledge can either be new or can interestingly temper, or confirm, knowledge elsewhere stated, but that may lack validation.

However, comparing PGx knowledge coming from distinct sources is challenging because of the heterogeneity of these sources. Indeed, such sources may use different vocabularies, different levels of granularity or even different languages to represent knowledge units. Therefore, there is a strong need for developing a common schema that would enable comparing knowledge extracted or discovered from various sources. Several

ontologies have already been developed for pharmacogenomics, but with different purposes, making them unadapted to the present need. In particular, SO-Pharm (Suggested Ontology for Pharmacogenomics) and PO (Pharmacogenomic Ontology) have been developed for knowledge discovery purposes rather than data integration or knowledge reconciliation [12, 13]. The PHARE ontology (for PHarmacogenomic Relationships) has been built for normalizing *gene – drug* and *gene – disease* relationships extracted from text and is not suitable for representing ternary PGx relationships [14]. More recently, Samwald et al. introduced the Pharmacogenomic Clinical Decision Support (or Genomic CDS) ontology, whose main goal is to propose consistent information about pharmacogenomic patient testing to the point of care, to guide physician decisions in clinical practice [15]. We have built PGxO by learning and adapting from these previous experiences. For consistency reasons and good practices, we mapped PGxO to concepts of these four pre-existing ontologies.

In this work, we propose to leverage Semantic Web and Linked Open Data (LOD) [1] technologies as a first step toward building a framework to represent and compare PGx relationships from various sources. We import knowledge of three origins to instantiate our “pivot” ontology, both illustrating the role of the ontology, and building a community resource for PGx research.

In the preliminary stage of this work [16], we proposed: (i) a first version of the PGxO ontology able to represent simple pharmacogenomic relationships and their potentially multiple provenances and (ii) a set of rules to reconcile PGx knowledge extracted from or discovered in various sources, i.e. to identify when two relationships refer to the same, or to different knowledge units. In this paper, we extend PGxO to improve its ability to represent PGx relationships extracted from the literature and by adding the notion of *proxy*. We experiment our approach by instantiating PGxO with knowledge of various provenances: PharmGKB, the biomedical literature, and results of studies that analyzed EHR data and linked DNA biobanks [11].

This paper is organized as follows. The “Methods” section introduces the methods used for the construction of PGxO, for encoding provenance information and for our rule-based approach to reconcile PGx knowledge. Details are also given about algorithms and techniques used to instantiate PGxO from the aforementioned sources. The “Results” section presents our ontology, PGxO, the reconciliation rules and results of the instantiation and reconciliation processes. Finally, we conclude this work by discussing the abilities of PGxO for representing and reconciling PGx knowledge and the next challenges to tackle.

Methods

Ontology construction

PGxO was manually and collaboratively developed by 3 persons (PM, CJ and AC) in 7 iterations (as of July 2018). We followed classical ontology construction methods and life cycle [17, 18], including the steps of specification, conception, diffusion and evaluation of the ontology.

Specification

Our aim is to represent and reconcile what we defined as PGx knowledge units. These are ternary relationships between one (or more) *genetic factor(s)*, one (or more) *drug treatment(s)* and one (or more) *phenotype(s)*. Such phenotypes can be the expected outcomes of the drug treatments or some adverse effects. In order to keep our ontology simple and leverage existing works representing PGx components, we restrain the *scope* of PGxO only to representing PGx knowledge units and not all facets of pharmacogenomics. The *objective* of PGxO is twofold: reconciling and tracing these PGx knowledge units. To enable this reconciliation, we need to encode metadata and provenance information about a PGx relationship.

Conception and diffusion

Because PGxO is of small size, the conception step was performed simultaneously with conceptualization, formalization and implementation steps. The ontology has been implemented in OWL using the Protégé ontology editor [19]. PGxO is conceived around the central class of `PharmacogenomicRelationship`, which enables associating two or three of the following key components of PGx: `Drug`, `GeneticFactor` and `Phenotype`. The expressive Description Logic (DL) associated with PGxO is $\mathcal{ALCI}(D)$ [20]. Successive versions of PGxO have been published online and shared with collaborators through both the NCBO BioPortal [21, 22] and GitHub [23]. We have followed [24] guidelines to report on the Minimum Information for the Reporting of an Ontology (MIRO) associated with PGxO and made this available at [23].

Evaluation

To evaluate our ontology, we used *competency questions* as proposed by Gangemi [25]. The questions we defined are the following:

- 1 Does PGxO enable to represent a PGx knowledge unit from the PGx state of the art (i.e. from a reference database or extracted from the biomedical literature), along with its provenance?
- 2 Does PGxO enable to represent a PGx knowledge unit discovered from clinical data, along with its provenance?
- 3 Does PGxO, coupled with its reconciliation rules, enable to decide if two knowledge units, with distinct provenances, may refer to the same thing?

We answered these questions twice, once early and once late in the iterations of the development of PGxO. For the former iteration, we manually instantiated PGxO with examples of knowledge units, associated with their provenances, from (i) PharmGKB, (ii) the literature (extracted by Semantic Medline [26] or FACTA+ [27]) and (iii) hand designed facts corresponding to what we thought may be discovered in EHRs. For the latter iteration, we answered these questions by instantiating PGxO with knowledge units extracted programmatically from PharmGKB and the biomedical literature, and manually from results reported by studies analyzing EHR data and linked biobanks. Details on the methods used to populate PGxO from these various sources are provided in following subsections.

Mappings

For consistency reasons and good practices, we manually mapped concepts of PGxO to the four aforementioned ontologies related to pharmacogenomics: SO-Pharm, PO, PHARE and Genomic CDS. These mappings are available in [28]. Because the NCBO BioPortal generates lexical-based mappings between the ontologies it hosts, it provides an initial set of mappings from PGxO to many standard ontologies. In particular, we manually completed PGxO BioPortal mappings to three standard and broad spectrum ontologies: MeSH, NCIt and SNOMED CT. These mappings are available in [29]. The two resulting sets of mappings are provided as independent OWL files to allow a flexible loading of the ontology with, or without mappings.

Provenance encoding

Data provenance (sometimes called lineage) refers to metadata that state where data came from, how they were derived, manipulated, and combined, and how they may have been updated [30]. With PGxO, we do not only want to represent units of knowledge of different origins, but also to trace their origins. Therefore, we need an encoding of the provenance of knowledge units. For this purpose, we leverage an existing ontology for provenance, named PROV-O [31], which is a W3C Recommendation since 2013. In addition, for some particular provenance metadata, PGxO reuses object properties of the high-level ontology DUL (Dolce+DnS Ultralight) [32].

PROV-O is built around three main concepts: `Entity`, `Activity` and `Agent`. Entities represent things that can be generated, modified, etc. by activities. Activities are realized by agents that can be either human or software agents. Entities can also be directly attributed to agents.

In terms of PROV-O concepts, authorities which publish sources from which we extract knowledge units are considered to be agents (e.g. the PharmGKB team, the National Library of Medicine (NLM) in charge of

PubMed, an hospital in charge of a repository of EHRs). Data sources (e.g. a version of PharmGKB, of PubMed, a repository of EHRs) are attributed to agents, and then may be used to derive data. These data, in turn, are used during the execution of an activity (e.g. a mining algorithm). Such execution generates entities that in our case are PGx knowledge units. Quantitative and qualitative metadata may be associated to an activity and to the entities it generates. For instance, one can specify the version of an algorithm, the date of its execution, the quality of the generated entities (such as their levels of confidence, their *p*-values or odds ratios). Thereby, a further comparison of two knowledge units having different or identical provenances may take into account these various elements.

Reconciliation rules

Besides representing and encoding heterogeneous PGx relationships within a single knowledge base, a comparison mechanism is needed to identify when two relationships refer to the same knowledge unit or not. However, the description of PGx relationships is highly heterogeneous depending on the sources we consider, leading to many relationships being similar to some extent, but not exactly identical. For example, one source may document a relationship between a gene variant gv_1 , a drug d_1 that causes a drug response phenotype p_1 , whereas a second source may only document the relationship between gv_1 and d_1 , omitting any drug response phenotype. Alternatively, a third source may document the same relationship at a broader level, for instance by mentioning only the involved gene g_1 , instead of stating the causative variant gv_1 (that is part of g_1).

To take into account this variability, we defined a set of rules enabling basic comparison mechanisms. The rules focus on identifying identical relationships, broader/narrower ones and related ones (to some extent). They compare two PGx relationships using their associated components (i.e. drugs, genetic factors, phenotypes). To represent and implement the defined rules, we investigated semantic web rule languages, such as SWRL and DL-Safe rules [33–35]. Unfortunately, we found them unadapted to our case, for expressiveness reasons. Indeed, those formalisms do not allow to check equalities or inclusions of sets of individuals, which instantiate defined DL expressions (see Additional file 1 for examples). Therefore, we represent the rules with our own formalism. On the left side of a rule, equalities or inclusions between sets of components of the two compared PGx relationships are tested and return a boolean result. Tests of equalities or inclusions can be combined using conjunctions (AND) and / or disjunctions (OR). On the right side, a link between the PGx relationships is to be added to the populated ontology if and only if the left side of the rule is

true. This link can specify the two PGx relationships as identical, related or one being broader than the other. As the rules are not formalized using a particular semantic web rule language, they are kept separated from the definitions of PGxO. Therefore, we implemented them in an independent Python program interacting with the populated ontology using the SPARQL query language. This program is executed once, derives novel facts from the rules and adds them to the knowledge base.

We take advantage of Semantic Web technologies, by using associated reasoning mechanisms for the comparison of PGx relationships. In particular, we use the semantics associated with `owl:sameAs` links that state that two URIs are actually referring to the same entity. The interpretation of the `rdfs:subClassOf` relation and its transitivity is also used when comparing a PGx relationship that may be more specific/general than another one.

Ontology instantiation

We instantiated the ontology with PGx knowledge units from various sources, first, to answer the *competency questions* defined for PGxO and, second, to build a data set called PGxLOD (*PGx Linked Open Data*) that we think may constitute a valuable community resource for PGx research.

With preexisting LOD

We initiated PGxLOD from a preexisting set of Linked Open Data made of interconnected genes, drugs, and diseases according to 6 standard databases [36]. Such LOD data set follows the Semantic Web standards, particularly by using the Resource Description Framework (RDF) language and Uniform Resource Identifiers (URI) [1], which makes it adequate to connect with Semantic Web ontologies.

This preexisting data set is an aggregation of data from ClinVar [37], DisGeNET [38], DrugBank [39], SIDER [40] and MediSpan (a proprietary database). Accordingly, it includes and relates data about drugs, diseases and phenotypes, but no PGx relationships. Nevertheless, this data set groups together data related to entities involved in PGx relationships, and mappings between entities that may be present in different data sources, e.g. a drug referenced both in DrugBank and SIDER. These mappings are of particular interest for our purpose of comparing PGx relationships, since those may be composed of entities referenced in these various sources.

The initial instantiation of PGxO with preexisting LOD is straightforward since entities representing genes, drugs and diseases are used to instantiate the corresponding PGxO concepts, using the RDF predicate `rdf:type`. In the following, we name PGxLOD v1 the result of this instantiation process. This constitutes the initial version

of PGxLOD, with no PGx relationships, to distinguish from PGxLOD v2, a version enriched with PGx relationships of various provenances.

With PharmGKB data

Second, PGxO was instantiated with data from PharmGKB [6], a reference database for pharmacogenomics. PharmGKB's *clinical annotations* describe PGx relationships between genes (potentially their variants), drugs, and phenotypes. They are produced by PharmGKB curators after a review of the biomedical literature and of recommendations of health agencies such as the US Food and Drug Administration. In addition, PharmGKB contains cross-references, i.e. identifiers of genes, variants, drugs and phenotypes within other databases (such as NCBI Gene for genes) or ontologies (such as the Anatomical Therapeutic Chemical Classification System for drugs).

Part of PharmGKB data are already available in the form of LOD as a part of the Bio2RDF project [41]. Nonetheless, this version is outdated and provides only a small portion of PharmGKB. Therefore, inspired by this precursor work and following the guidelines of the Bio2RDF project, we developed new scripts producing a more complete RDF version of PharmGKB. These scripts transform the latest downloadable text files of PharmGKB, first, into a SQL database (with a script named *pharmgkb2sql*) and then into RDF triples (with a script named *pharmgkb2triples*).

Drug response phenotypes provided in clinical annotations by PharmGKB are non easy to translate as they are reported within plain-text sentences. Because PharmGKB also provides a broad type of drug responses (Efficacy, Toxicity/ADR, Metabolism/PK, Dosage, Other) in a structured manner, we decided, for simplicity, to use those directly instead of further text mining analysis on textual descriptions, and then considered only Efficacy and Toxicity/ADR, because they encompass the drug response phenotypes we want to capture. Accordingly, PGx relationships in PGxLOD are associated with these two high level types of drug responses.

Components of PGx relationships (i.e. drugs, genes and variants) are represented with URIs using both PharmGKB identifiers and Bio2RDF naming conventions. In addition, PharmGKB cross-references to external databases and ontologies are used to map PharmGKB URIs either to URIs already defined in our LOD, or to new ones. The type of relation used is `bio2rdf:x-ref` for every cross reference; doubled with either a `owl:sameAs` relation when the cross-reference points to an identical entity in an external database, or with a `rdf:type` relation when it points to an ontology class.

Among the metadata associated with PharmGKB clinical annotations, we particularly keep the *level of evidence*. Levels of evidence have been defined by PharmGKB [6] as

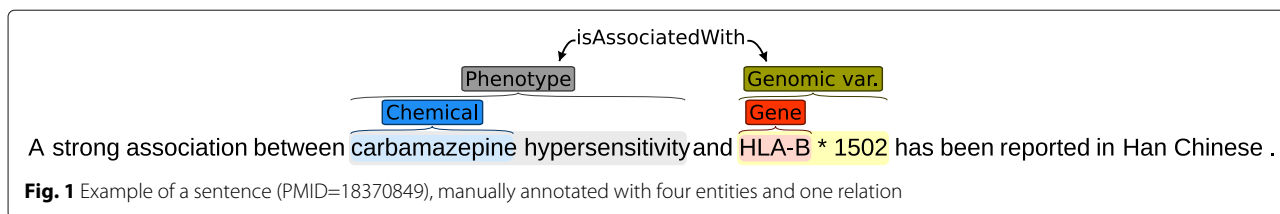
a six-level scale (1A, 1B, 2A, 2B, 3, 4), where higher levels (1A, 1B, 2A, 2B) indicate that a relationship has been significantly studied or is medically implemented; and lower levels (3, 4) indicate that the PGx relationship has only been reported in a single study or lacks clear evidence. Levels of evidence are of particular importance as they may help us identify PGx knowledge with irregular levels of validation in various sources.

With the biomedical literature

Third, we instantiated PGxO with elements extracted automatically from the biomedical literature. Here we used a supervised machine learning prototype for relation extraction from text, trained on a manually annotated corpus. Please note that in this work, we are not aiming at achieving the best performance, but rather aiming at showing that we can extract PGx relationships from text, along with their provenance metadata, and compare these to others, extracted from distinct sources. This illustrates that PGxO enables structuring, documenting (the name of the algorithm used, its performance, etc.), then comparing a relationship extracted from text.

We assembled a set of 657,538 sentences from 86,520 PubMed abstracts related to pharmacogenomics. Removing malformed sentences, based on tokenization errors, and sentences that do not contain at least one drug and one genetic factor, based on named entities recognized by PubTator [42], we obtained a corpus of 176,704 sentences. Out of those, we randomly selected 307 sentences and had each sentence manually annotated with the BRAT software [43], by 3 distinct annotators from a group of 11 pharmacists, biologists and bioinformaticians. The annotation task is precisely specified in annotation guidelines, publicly available [44]. Annotations are limited here to four broad entity types, mainly involved in PGx relationships: Genes, Genomic Variations, Drugs and Phenotypes and to two broad relation types, "isAssociatedWith" and "isEquivalentTo", between these entities. The latter is used only to relate the numerous acronyms to their extended form. An example of an annotated sentence is shown in Fig. 1, and the main characteristics of the corpus are summarized in Table 1.

Our approach is classically composed of two steps: a Named Entity Recognition (NER), followed by a relation extraction. Two supervised machine learning models were trained for the first step, and a third one for the second step. The NER is performed using a Convolutional Neural Network (CNN) model described in [45], trained on the 307 annotated sentences. We first annotate shallow entities (Fig. 1) using an instance of this model with *PubTator* annotations as an input. We name these entities *first - layer* entities. Then, a second instance of the same model is trained to annotate *second - layer* entities, i.e. entities with an offset that includes a *first - layer*



entity, using the *first – layer* entities as input. In Fig. 1, *first–layer* entities would be *carbamazepine* and *HLA–B*, and *second – layer* entities are *carbamazepine hypersensitivity* and *HLA-B*1502*. Finally, we trained a model similar to the one described in [46] to extract relationships between identified entities.

Reasonable meta-parameters were selected according to previous experiments. The size of the word embeddings was set to 100. The size of the *PubTator* and *first – layer* embeddings was set to 20. The kernel size of the convolution was set to 100. Word embeddings were pre-trained on ~3.4 million PubMed abstracts (corresponding to all those published between Jan. 1, 2014 and Dec. 31, 2016) using the method described in [47]. Performances of the two steps, evaluated on a 10-fold cross validation, are respectively reported in Table 2 and 3. In these two tables, the f1-score is defined as the harmonic mean of the precision and recall, i.e. $f1\text{-score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$.

Trained models are applied to a test set of 176,704 sentences of PubMed abstracts, to extract automatically relations from text. After filtering out relationships that relate two GenomicVariations, two Phenotypes, or two Drugs, both the manually annotated relations and the automatically extracted ones are transformed into RDF.

Types of entities are manually aligned to the corresponding classes of PGxO. Each annotated and extracted entity is associated with an URI that is constructed, depending on its type, either from an identifier of a reference database (such as NCBI Gene for genes) or from an identifier of an ontology (such as ATC for drugs). Distinct reference databases or ontologies may be used for each

type of entities depending on their availability. Accordingly, we defined an arbitrary order of choice for searching for references, presented in Table 4. This choice was motivated in part by the output of *PubTator*. For each type, the procedure is the following: given an entity, the first reference database or ontology is searched for the entry using string matching; if no entry matches, the next reference database or ontology is searched. Lastly, if no entry is found, we create a local URI within the PGxLOD namespace. Considering an extracted entity, when an entry is found in a reference database, its identifier is used to construct the corresponding URI. When an entry is found in an ontology (i.e. a class of an ontology), the extracted entity is given a local URI, and instantiates the ontology class.

With electronic health records and linked biobanks studies

Fourth, we instantiated PGxO with PGx knowledge manually extracted from the reading of ten studies on patient Electronic Health Records (EHRs) and linked biobanks [9, 11, 48–55]. For instance, Kawai et al. [53] report a statistical association (OR=2.05, 95%) between the haplotype CYP2C9*3, and severe bleeding, in patients treated with warfarin. Their study was performed on a biobank named BioVU, linked to patient EHRs of the Vanderbilt University Hospital [56]. The ten studies were selected from mentions in CPIC (Clinical Pharmacogenetics Implementation Consortium) guidelines [3] and in the literature review of Denny et al. [57].

Entities involved in PGx relationships were manually associated with URIs already defined in PGxO and PGxLOD. The aim here is to assess the adequacy of PGxO to represent results of such studies. Indeed, it is

Table 1 Statistics of named entities and relations manually annotated in our 307-sentence corpus

Named entities			Relations	
Type	First-layer	Second-layer	Type	
Gene	452	20	isAssociatedWith	582
GenomicVariation	74	166	isEquivalentTo	77
Drug	459	36		
Phenotype	262	251		
Total		1720	Total	659

Entities annotated in multiple sentences are counted multiple times. Second-layer entities refer to entities which offset includes the annotation of a first-layer entity

Table 2 Named Entity Recognition (NER) performance in terms of precision (P), recall (R) and f1-score (F1)

	P	R	F1	std
Drug	0.92	0.87	0.89	0.03
Gene	0.97	0.91	0.94	0.03
Phenotype	0.84	0.66	0.74	0.09
Genomic variation	0.81	0.69	0.74	0.08
All entities	0.86	0.80	0.83	0.05

Results of second-layer entities take into account the prediction error of the first-layer entities. Std stands for F1 standard deviation

Table 3 Relation extraction performance in terms of precision (P), recall (R) and f1-score (F1)

	P	R	F1	std
isAssociatedWith	0.61	0.35	0.44	0.08
isEquivalentTo	0.73	0.78	0.75	0.14
All relations	0.67	0.56	0.61	0.08

Std stands for F1 standard deviation. Results take into account the prediction error for the entities

one of our use cases to consider PGx researchers who want to compare the results they obtained on their local biobanks+EHRs, to results elsewhere reported.

Results

PGxO

Illustrated in Fig. 2, PGxO is composed of eleven concepts (owl:Class), organized around the central concept PharmacogenomicRelationship, which represents an atomic unit of PGx knowledge.

Instances of the concepts may be related by various types of relations (i.e. owl:ObjectProperty). Relation types causes and isCausedBy are used to connect components of a PGx relationship and are defined as inverses: causes ≡ isCausedBy⁻. The relation type partOf (or ro:BFO_0000050), from the Relation Ontology (RO) [58], is used to express that a genomic variation is located within the sequence of a specific gene. The relation type dependsOn (ro:RO_0002502), also from RO, is used to express complex phenotypes that involve other entities, e.g. gene expression such as *the expression of VKORC1* or drug response phenotypes such as *carbamazepine hypersensitivity*.

The concept PharmacogenomicRelationship is described by a Description Logics [20] axiom that may be decomposed into a union of three other axioms. Indeed, we consider that a PGx relationship involves at least two or three of the following components: drugs, genetic factors and phenotypes (i.e. drug response phenotypes). For more flexibility, we allow drug components to be either a drug, or something that depends on a drug (Axiom 1). Similarly, we allow genetic factor components to be either a genetic factor, or something that depends on a genetic

Table 4 Reference databases and ontologies used to normalize the entities extracted from text

Order	Drug	Gene	GenomicVariation	Phenotype
1 st	MeSH	NCBI Gene	dbSNP	MeSH
2 nd	ChEBI	PGxLOD	PGxLOD	MEDDRA
3 rd	ATC			PGxLOD
4 th	PGxLOD			

PGxLOD means that a local URI is created

factor (Axiom 2). This flexibility allows to represent relationships that involve, for instance, the expression of a gene (something that depends on a genetic factor, e.g. *the expression of VKORC1*) instead of a gene, or a drug resistance or sensitivity (something that depends on a drug, e.g. *carbamazepine hypersensitivity*) instead of a drug.

Axiom 1 $DComponent \equiv Drug \sqcup \exists dependsOn.Drug$

Axiom 2 $GFComponent \equiv GeneticFactor \sqcup \exists dependsOn.GeneticFactor$

Another flexibility resides in allowing a PGx relationship to have only two of these three components. Indeed, one component may be missing for multiple reasons: the relationship may still be under study and some of its components unknown; we can also expect errors during automatic extraction leading to missing components. Therefore, Axioms 3, 4 and 5 represent the three possibilities of having two components among the three considered. Axiom 6 is the final axiom describing the concept PharmacogenomicRelationship. A PGx relationship involving the three types of components trivially validates the axioms.

Axiom 3 $PR_1 \equiv \exists causes.Phenotype \sqcap \exists isCausedBy.DComponent$

Axiom 4 $PR_2 \equiv \exists causes.Phenotype \sqcap \exists isCausedBy.GFComponent$

Axiom 5 $PR_3 \equiv \exists isCausedBy.DComponent \sqcap \exists isCausedBy.GFComponent$

Axiom 6 $PharmacogenomicRelationship \sqsubseteq PR_1 \sqcup PR_2 \sqcup PR_3$

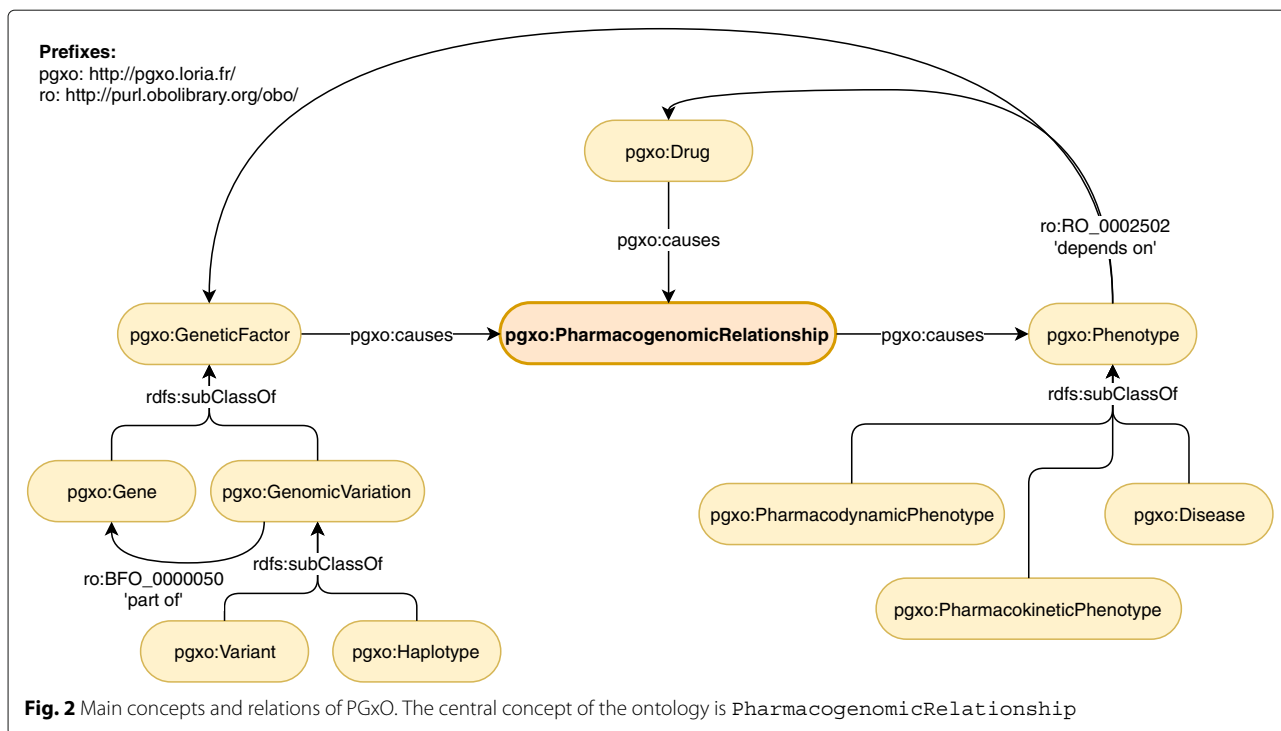
Axiom 6 is not a concept definition (using ≡) because we consider that presenting two (or three) components is a required condition, but not a sufficient one to be a PharmacogenomicRelationship. Examples of encoding of PGx relationships and their provenance are detailed in the next subsection.

PGxLOD: an instantiation of PGxO with knowledge units of various provenances

PGxLOD is the knowledge base that results from the instantiation processes of PGxO. Because a large portion of the data of PGxLOD has no license restriction, we provide an open access to this portion (only) at <https://pgxlod.loria.fr>. Full access to PGxLOD that includes the small set of proprietary data is granted upon request. The population processes of PGxLOD are detailed in the next paragraphs and their results are summarized in Table 5.

With preexisting LOD

Table 6 summarizes results of the instantiation of PGxO with our preexisting LOD. At this stage PGxLOD does not contain any PGx relationship, but provides entities appearing as components of PGx relationships, as well as mappings between these entities.



With PharmGKB data

Table 7 summarizes results of the instantiation of PGxO with PharmGKB data (2018-03-05 release).

Figure 3 provides an example of a PharmGKB relationship represented with PGxO. It represents a relationship between the haplotype TPMT*3C and the drug azathioprine. This relationship was extracted with our algorithm named *pharmgkbsql2triples*. Accordingly this algorithm

is specified as provenance metadata, along with its version and the version of PharmGKB. This allows coexisting extractions of concurrent versions of PharmGKB or the algorithm. The level of evidence of the relationship is represented with PROV-O concepts.

With the biomedical literature

We instantiated PGxO with the manually annotated sentences of our 307-sentence corpus, and with the output of our prototype for relation extraction on a test corpus formed by the all 176,704 sentences unannotated or annotated. We extracted from these sentences, 51,924 entities (8412 drugs, 10,812 genes, 8740 genomic variations and 23,960 phenotypes) and 65,182 PGx relationships between them. Table 8 shows the statistics for the

Table 5 Main statistics of PGxLOD v2

PGxO concept	Number of instances
Drug	51,459
GeneticFactor	386,801
Gene	172,881
GenomicVariation	213,910
Haplotype	33
Variant	204,875
Phenotype	88,247
Disease	47,573
PharmacodynamicPhenotype	63
PharmacokineticPhenotype	44
PharmacogenomicRelationship	68,431
from PharmGKB	2701
from the literature	65,720
from EHR studies	10

Table 6 Statistics of the instantiation of PGxO with data from PGxLOD v1

Source	Genes	Variants	Drugs	Diseases	Phenotypes
ClinVar	21,487	103,219	0	0	6837
DisGeNET	85,893	49,279	0	38,727	6092
DrugBank	4300	0	7740	0	0
MediSpan	0	0	5820	2481	0
SIDER	0	0	25,479	6291	0
UniProt	25,456	0	0	0	0
Total	137,136	152,498	39,039	47,499	12,929

Table 7 Numbers of PGx relationships extracted from PharmGKB v2018-03-05

	Caused phenotype		Level of evidence					All	
	Toxicity/ADR	Efficacy	1A	1B	2A	2B	3		4
# PGx relationships	1268	1531	44	11	71	97	2270	208	2701

Some PGx relationships can cause both *Toxicity/ADR* and *Efficacy*

normalization of these entities to identifiers of reference databases or ontologies listed in Table 4. Figure 4 illustrates the RDF encoding of a PGx relationship extracted from the literature. It is noteworthy that the type of relation is encoded in the provenance metadata. Our prototype only outputs relations of the broad type “isAssociatedWith”, but other types could be expected with a more

advanced system, e.g. “increases” or “decreases”. Figure 4 also illustrates how the entity representing the TPMT gene reuses an URI from the Bio2RDF transformation of the NCBI Gene database, while the entity representing the 6-mercaptopurine instantiates a MeSH class. This differentiates the use of reference databases or ontologies when normalizing.

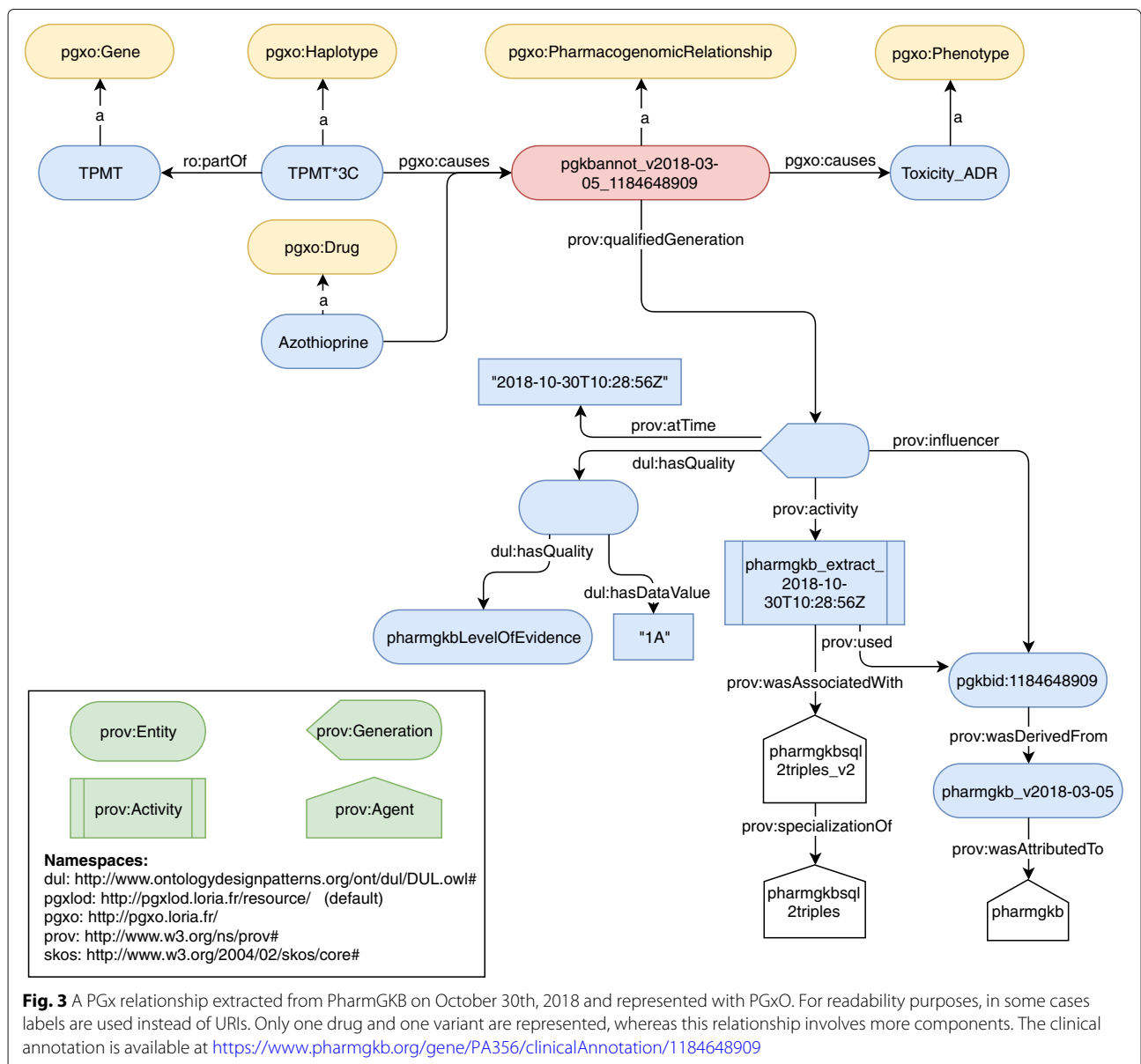


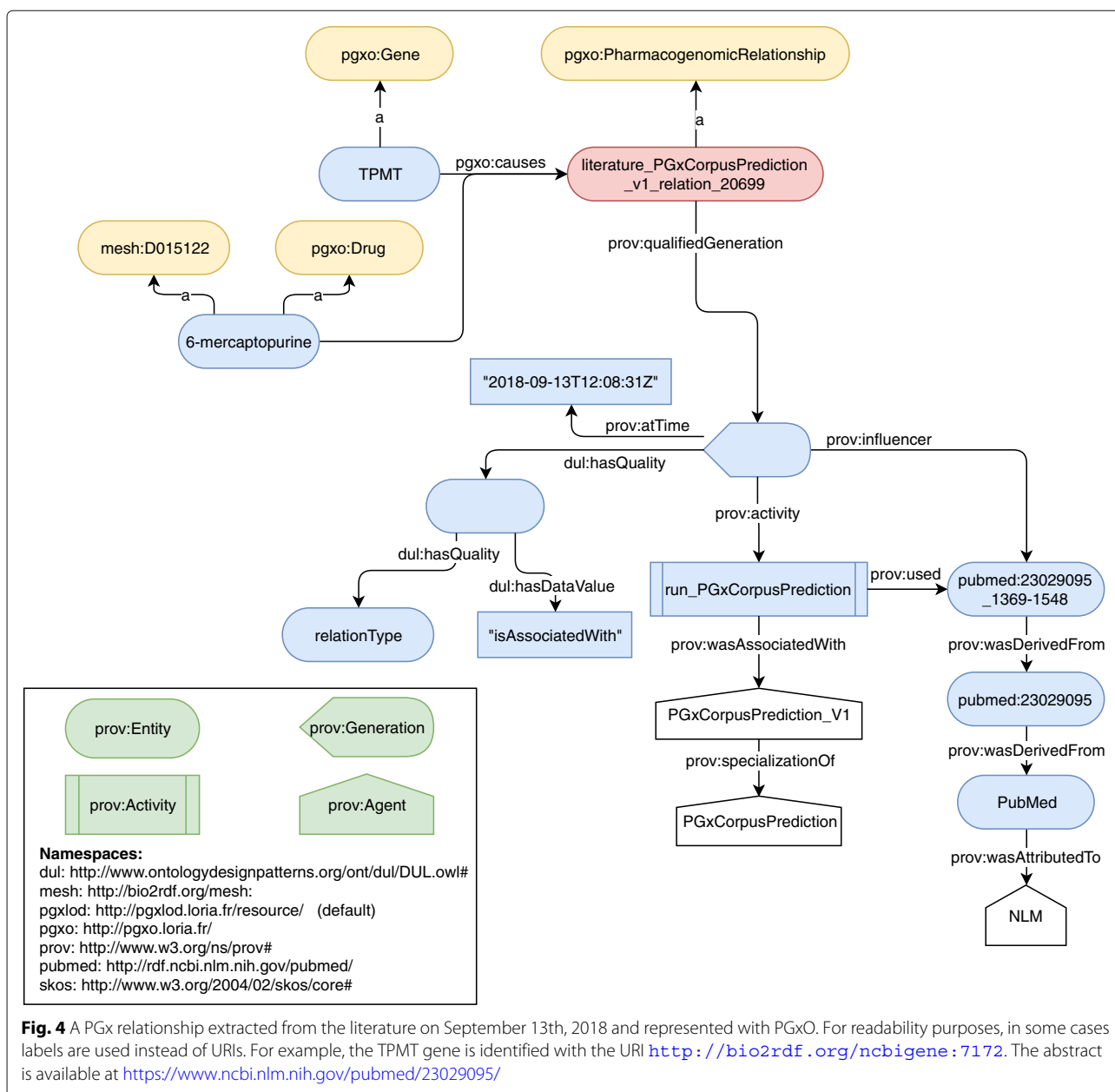
Table 8 Numbers of unique entities recognized in the test corpus and successfully mapped with reference databases or ontologies

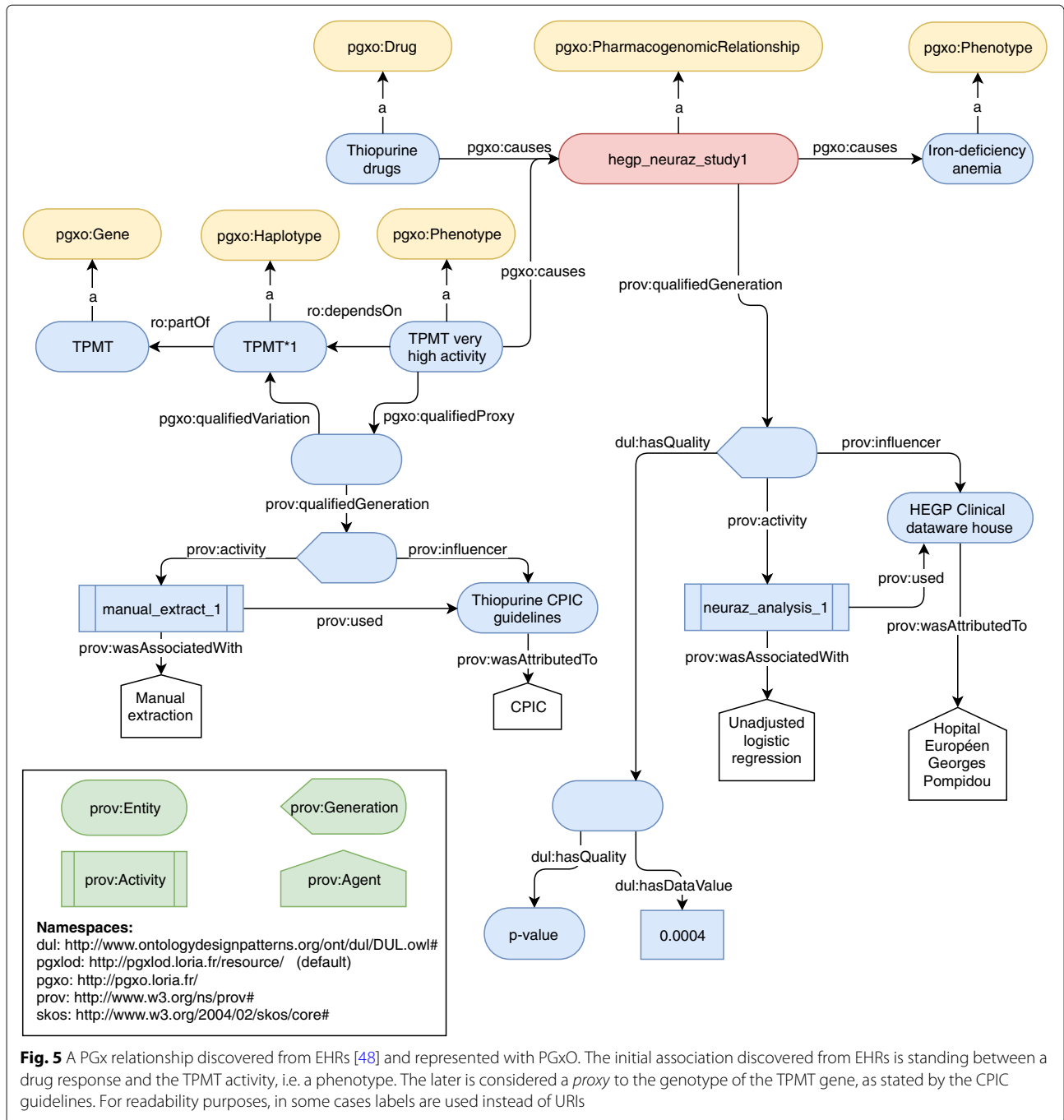
Database / Ontology	Drug	Gene	GenomicVariation	Phenotype
MeSH	1600	n/a	n/a	1625
ChEBI	285	n/a	n/a	n/a
ATC	78	n/a	n/a	n/a
NCBI Gene	n/a	4907	n/a	n/a
dbSNP	n/a	n/a	803	n/a
MEDDRA	n/a	n/a	n/a	0
PGxLOD	6449	5905	7937	22,335
Total	8412	10,812	8740	23,960

Reference databases and ontologies are listed in Table 4

With electronic health records and linked biobank studies

Each of the ten studies listed previously results in one instance of a PGx relationship, along with its provenance. Interestingly, out of ten, eight were derived from the BioUV biobank and its linked EHRs [56], one from clinical data of the eMERGE Network [59] and one from data of the HEGP, a French University Hospital [60]. Out of the same ten relationships, six were obtained from a statistical analysis using linear regression and four using logistic regression. Regarding genetic factors, relationships involve either a single nucleotide polymorphism (7/10), an haplotype (2/10) or an enzyme activity (1/10). For example, Fig. 5 represents the instantiation of PGxO, achieved from the results of Neuraz et al. [48] and the thiopurine





CPIC guidelines. In this particular case, no genetic data was provided in the study, but an enzyme activity. Indeed the TPMT enzyme activity may be considered as a *proxy* for the genotype of the TPMT gene, as stated in the thiopurine-related CPIC guidelines [61]. We added a RDF triple stating that the TPMT activity depends on the TPMT haplotype (with the `ro:dependsOn` relation type), and documented the provenance of this

assertion (with the `pgxo:qualifiedProxy` and `pgxo:qualifiedVariation` relation types and PROV-O concepts and relation types). This representation is possible because our axioms defining a PGx relationship allow using a phenotype as a proxy for a genetic factor instead of a genetic factor itself. We expect this to facilitate later comparison of the results of studies without genetic data, to the state of the art.

Reconciliation rules

Definition and implementation of the reconciliation rules

We defined five rules for simple pair-wise comparison of PGx relationships. These rules are able to identify when two PGx relationships with distinct provenances are in fact referring to the same knowledge unit, to a more specific knowledge unit, or to related knowledge units (to some extent). Indeed, among the five rules, Rule 1 is dedicated to identifying identical relationships, Rules 2, 3, and 4 to identifying broader/narrower ones and Rule 5 to identifying PGx relationships related by some of their components. All five rules are provided in Additional file 1 as well as examples of their application on RDF graphs. In the next paragraphs, as an example, the simpler rule, Rule 1, identifying when two PGx relationships are referring to the same knowledge unit is presented and illustrated. Other rules are a bit more complex, but follow the same principles.

Rules compare PGx relationships on the basis of their components, i.e. sets of drugs, genetic factors and phenotypes. Accordingly, considering r , an instance of the PharmacogenomicRelationship concept from a knowledge base \mathcal{KB} , the following sets are defined.

Notation 1 We denote D , the set of instances of Drug that cause r , defined as $D = \{d \mid \mathcal{KB} \models \text{Drug}(d) \text{ and } \mathcal{KB} \models \text{causes}(d, r)\}$

Notation 2 We denote G , the set of instances of GeneticFactor that cause r , defined as $G = \{g \mid \mathcal{KB} \models \text{GeneticFactor}(g) \text{ and } \mathcal{KB} \models \text{causes}(g, r)\}$

Notation 3 We denote P , the set of instances of Phenotype caused by r , defined as $P = \{p \mid \mathcal{KB} \models \text{Phenotype}(p) \text{ and } \mathcal{KB} \models \text{causes}(r, p)\}$

Therefore, when comparing two PGx relationships denoted by r_1 and r_2 , the sets of their components will be denoted D_1, G_1, P_1 and D_2, G_2, P_2 . The first reconciliation

rule identifies when two PGx relationships are referring to the same knowledge unit; it is defined as follows:

Rule 1 $D_1 = D_2 \text{ AND } G_1 = G_2 \text{ AND } P_1 = P_2 \Rightarrow \text{owl:sameAs}(r_1, r_2)$

This rule states that when two relationships involve the same sets of drugs, of genetic factors and of phenotypes, they refer to the same knowledge unit. Therefore, the link $\text{owl:sameAs}(r_1, r_2)$ should be added to the knowledge base. For example, consider the RDF graph presented in Fig. 6. We have:

- $D_1 = D_2 = \{\text{warfarin}\}$
- $G_1 = G_2 = \{\text{CYP2C9}\}$
- $P_1 = P_2 = \{\text{cardiovascular_diseases_inst1}\}$

Therefore, the left part of Rule 1 is true, and the link $\text{owl:sameAs}(r_1, r_2)$ should be added to the knowledge base.

Other rules, details and examples are available in Additional file 1. Rules 2, 3 and 4 conclude in indicating that a relationship is more specific than the other by adding the link $\text{skos:broadMatch}(r_1, r_2)$ to the knowledge base. Rule 5 concludes that they are related by adding the link $\text{skos:relatedMatch}(r_1, r_2)$.

Execution of the reconciliation rules on PGxLOD

We executed our reconciliation rules on PGxLOD v2 (Table 5). As each relationship is compared to all the others, this led to $68,430 \times 68,431 = 4,682,733,330$ comparisons performed.

This execution generated owl:sameAs links (Table 9), skos:broadMatch links (Table 10) and skos:relatedMatch links between the PGx relationships in PGxLOD. We observed skos:relatedMatch (37,758) only between pairs of relationships extracted from the

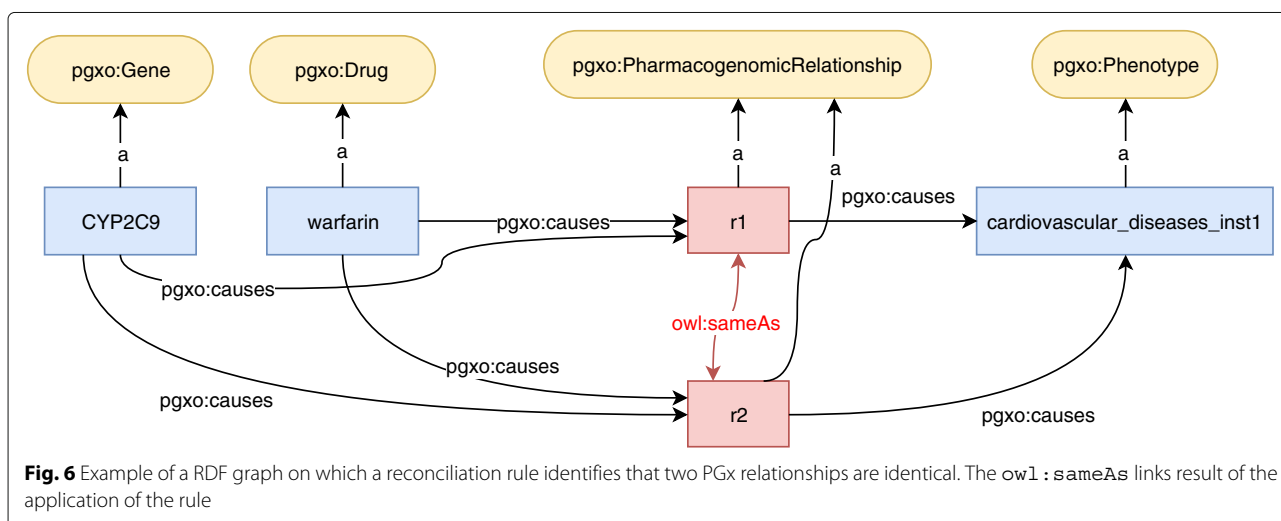


Table 9 Number of owl:sameAs links between PGx relationships from each source

	EHRs	Literature	PharmGKB
EHRs	0	0	0
Literature	0	109,078	0
PharmGKB	0	0	132

literature. Interestingly, for 132 PGx relationships from PharmGKB an identical relationship was found (generating 132 owl:sameAs links). Additionally, 14 sentences from the biomedical literature are identified as more generic than what is stated in the EHR+biobank studies.

We can notice that skos:broadMatch links exist between different sources while owl:sameAs and skos:relatedMatch links only refer to PGx relationships from the same sources. Some possible explanations reside in a lack of mappings between entities in different vocabularies and the use of broad phenotypes for PGx relationships from PharmGKB (i.e. Toxicity/ADR, Efficacy) that are distinct from more specific phenotypes elsewhere stated such as, for example, cardiovascular diseases.

Discussion

Instantiating PGxO with knowledge extracted from various sources allows to answer the defined competency questions: we are able to represent PGx relationships extracted either from the state of the art (reference databases or the literature) as well as from EHR+biobank studies. The use of heterogeneous sources for instantiating our ontology improved in several ways the modeling of PGx relationships previously drafted in [16]. Among other things, we enabled the representation of phenotypes as proxies for a specific genotype, such as an enzyme activity. The encoding of metadata has also been enriched to enable the encoding of the various metrics associated with the different kinds of knowledge extractions.

By using Semantic Web technologies, our global framework for knowledge comparison in PGx can easily leverage knowledge defined elsewhere such as ontologies or other available LOD sets. This is of particular importance

Table 10 Number of skos:broadMatch links between PGx relationships from each source

	EHRs	Literature	PharmGKB
EHRs	0	14	0
Literature	0	133,966	0
PharmGKB	0	865	98

Rows represent origins of the links and columns represent destinations

as the reconciliation rules depend on existing mappings and subsumption relations. Moreover, the proposed encoding can easily evolve depending on one’s needs. However, in a data warehousing perspective, Semantic Web technologies require high data maintenance to follow the evolution of associated databases, LOD sets and ontologies. Therefore, one challenge is to keep PGxLOD up-to-date with respect to the associated data sources.

Several directions are considered to continue this work. Regarding the extraction from PharmGKB, more detailed drug response phenotypes could be extracted from the plain-text sentences describing the clinical annotations in the database. This would require text mining similar to what we have done for processing literature but this would enable a more accurate comparison between the content of PharmGKB and other sources.

Our prototype for knowledge extraction from the literature constitutes solely a baseline. It faces limitations relatively easy to improve. First, the NER model, in its current form, does not detect discontinuous entities that may appear in the literature (such as “the response of the selective serotonin reuptake inhibitors paroxetine” where the entity “response of paroxetine” is discontinuous). This is a limitation since missed entities lead to missed relationships. In addition, the two steps NER procedure can only detect fairly simple included entities. In practice, multiple levels of inclusion can be observed in the literature and cannot be captured by our system. Finally, a larger training corpus would improve the performance of the learned models, since deep learning architectures usually require large annotated corpora in order to achieve reasonable performances. In addition, the normalization process, which results are reported in Table 8, is naive and may be improved using lexical resources and ontology repositories. Those list term synonyms, variants and bridges between these resources.

The manual instantiation of PGxO with knowledge extracted from EHRs constitutes only a proof of concept. One notable drawback is that gene variants and precise drug response phenotypes are not available in most cases. Thus, the knowledge discovery process needs to rely on proxies such as a phenotype being a marker of the patient genotype or a stable dose requirement being a marker of the patient sensitivity to the considered drug. Therefore, a PGx relationship discovery from EHRs would benefit from a more complete list of proxies. To our knowledge, no such list is available. In addition, more contextual information about knowledge discovered from patient data would be of interest. For example, the ethnicity of patients [53] or the indications for which patients are treated [9] may be necessary to properly document some PGx relationships. Considering these challenges, one perspective of the current work relies in automatically instantiating PGxO with knowledge extracted by mining EHRs.

The proposed reconciliation rules were executed on PGxLOD, providing first results of reconciliation. However, to compare PGx relationships involving entities from different vocabularies, the rules rely on the existence of equivalence or subsumption relationships between the URIs of these entities. Therefore, a major task and perspective resides in completing the mappings between entities of various provenances. Using both concept hierarchies and ontology-to-ontology mappings defined in the UMLS [62] or the NCBO Bioportal [63] may improve knowledge comparison. Especially, this may be particularly useful when considering knowledge extracted from EHRs, which are expressed with concepts of ontologies used in the encoding of clinical practice such as ICD or RxNorm. Finally, the reconciliation rules strictly compare the components of a PGx relationships: drugs, genetic factors, and phenotypes. However, other features could be considered, such as the specific chemicals of a drug. Such features could be involved in a fuzzy comparison highlighting similar but not strictly equivalent relationships.

Conclusions

In this article, we presented a simple ontology called PGxO to represent pharmacogenomic knowledge and its provenance. With the combined use of PROV-O and DUL, we demonstrated that PGxO can structure knowledge extracted from various sources such as reference databases (i.e. PharmGKB), the literature, clinical guidelines or EHR+biobank studies. We also defined and implemented a set of rules allowing to compare and reconcile PGx knowledge units from different sources. PGxO and the reconciliation rules constitute a first step to a semantic framework able to represent, trace, confront and reconcile PGx relationships from various origins. A first experiment with these rules highlights equivalent and comparable pieces of knowledge across various data sources, opening perspectives for fine grained comparison and interpretation of the content of PGx sources. Finally, we think that the resulting and integrated data set called PGxLOD constitutes by itself a valuable resource for PGx research. This data set is made available to the community and will be improved with additional knowledge from the state of the art and from EHR mining.

Additional file

Additional file 1: PGxO reconciliation rules. This PDF file provides the definition of the five reconciliation rules and illustrates their behavior with concrete examples. (PDF 337 kb)

Abbreviations

ADR: Adverse drug reaction; CNN: Convolutional neural network; CPIC: Clinical pharmacogenetics implementation consortium; DL: Description logic; DNA: Deoxyribonucleic acid; EHR: Electronic health record; F1: F1-score; LOD: Linked open data; NCBI: National center for biotechnology information; NCBO:

National center for biomedical ontology; NER: Named entity recognition; NLM: National library of medicine; OWL: Web ontology language; P: Precision; PGx: Pharmacogenomics; PGxLOD: PGx linked open data; PGxO: PGx ontology; PHARE: PHarmacogenomic Relationships ontology; PharmGKB: The pharmacogenomics knowledge base; PROV-O: The PROV ontology; R: Recall; RDF: Resource description framework; RO: Relation ontology; SPARQL: SPARQL protocol and RDF query language; SWRL: Semantic web rule language; URI: Uniform resource identifier; UMLS: Unified medical language system; W3C: World wide web consortium

Acknowledgements

We acknowledge the participants of the NETTAB 2017 workshop for their constructive feedback on the preliminary results of this work. This work was conducted using the Protégé resource, which is supported by grant GM10331601 from the National Institute of General Medical Sciences of the United States National Institutes of Health.

Funding

Publication charges are supported by Inria DGDS-IES. This work is supported by the *PractiKPharma* project, founded by the French National Research Agency (ANR) under Grant No. ANR-15-CE23-0028, by the IDEX "Lorraine Université d'Excellence" (15-IDEX-0004), by the *Snowball* Inria Associate Team, and by the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No. 701771.

Availability of data and materials

Our ontology, PGxO, is available on the NCBO BioPortal [22] and GitHub [23]. Our set of PGx linked data, PGxLOD, is available at <https://pgxlod.ioria.fr/>.

About this supplement

This article has been published as part of *BMC Bioinformatics Volume 20 Supplement 4, 2019: Methods, tools and platforms for Personalized Medicine in the Big Data Era (NETTAB 2017)*. The full contents of the supplement are available online at <https://bmcbioinformatics.biomedcentral.com/articles/supplements/volume-20-supplement-4>.

Authors' contributions

PM, AN and AC designed the experiments. PM, JL and AC wrote the manuscript. PM, CJ and AC designed the ontology and the encoding of the metadata. PM, AN and AC designed the reconciliation rules. PM implemented and ran the process of knowledge reconciliation. PM and GH implemented and ran the processes of knowledge extraction from PharmGKB. JL implemented and ran the processes of knowledge extraction from the biomedical literature. AC extracted knowledge related to biobanks and EHRs studies. PR set up and administers the virtuoso server for PGxLOD. AT and CJ advised on the definition of mappings between LOD entities. All authors read and approved the final manuscript.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Université de Lorraine, CNRS, Inria, LORIA, 54000 Nancy, France. ²LIRMM, Université de Montpellier, CNRS, 34095 Montpellier, France. ³Stanford Center for Biomedical Informatics Research, Stanford University, 94305 Stanford, California, USA.

Published: 18 April 2019

References

- Bizer C, Heath T, Berners-Lee T. Linked data-the story so far. *Int J Semant Web Inf Syst.* 2009;5(3):1–22.
- Xie H-G, Frueh FW. Pharmacogenomics steps toward personalized medicine. *Personalized Med.* 2005;2(4):325–37.
- Caudle KE, Klein TE, Hoffman JM, Muller DJ, Whirl-Carrillo M, Gong L, et al. Incorporation of pharmacogenomics into routine clinical practice: the Clinical Pharmacogenetics Implementation Consortium (CPIC) guideline development process. *Curr Drug Metab.* 2014;15(2):209–17.
- Martin MA, Hoffman JM, Freimuth RR, Klein TE, Dong BJ, Pirmohamed M, et al. Clinical Pharmacogenetics Implementation Consortium Guidelines for HLA-B Genotype and Abacavir Dosing: 2014 update. *Clin Pharmacol Ther.* 2014;95(5):499–500.
- Amstutz U, Henricks LM, Offer SM, Barbarino J, Schellens JHM, Swen JJ, et al. Clinical Pharmacogenetics Implementation Consortium (CPIC) Guideline for Dihydropyrimidine Dehydrogenase Genotype and Fluoropyrimidine Dosing: 2017 Update. *Clin Pharmacol Ther.* 2018;103(2):210–6.
- Whirl-Carrillo M, McDonagh EM, Hebert JM, Gong L, Sangkuhl K, Thorn CF, et al. Pharmacogenomics knowledge for personalized medicine. *Clin Pharmacol Ther.* 2012;92(4):414.
- Garten Y, Coulet A, Altman RB. Recent progress in automatically extracting information from the pharmacogenomic literature. *Pharmacogenomics.* 2010;11(10):1467–89.
- Ioannidis JP. To replicate or not to replicate: the case of pharmacogenetic studies: Have pharmacogenomics failed, or do they just need larger-scale evidence and more replication? *Circ Cardiovasc Genet.* 2013;6(4):413–8.
- Delaney JT, Ramirez AH, Bowton E, Pulley JM, Basford MA, Schildcrout JS, et al. Predicting clopidogrel response using DNA samples linked to an electronic health record. *Clin Pharmacol Ther.* 2012;91(2):257–63.
- Ramirez AH, Shi Y, Schildcrout JS, Delaney JT, Xu H, Oetjens MT, et al. Predicting warfarin dosage in European-Americans and African-Americans using DNA samples linked to an electronic health record. *Pharmacogenomics.* 2012;13(4):407–18.
- Birdwell KA, Grady B, Choi L, Xu H, Bian A, Denny JC, et al. The use of a DNA biobank linked to electronic medical records to characterize pharmacogenomic predictors of tacrolimus dose requirement in kidney transplant recipients. *Pharmacogenet Genomics.* 2012;22(1):32–42.
- Coulet A, Smail-Tabbone M, Napoli A, Devignes M-D. Suggested Ontology for Pharmacogenomics (SO-Pharm): Modular Construction and Preliminary Testing. In: *On the Move to Meaningful Internet Systems 2006: OTM 2006 Workshops*, Montpellier, France, October 29 - November 3, 2006. Proceedings, Part I. Springer; 2006. p. 648–57.
- Dumontier M, Villanueva-Rosales N. Towards pharmacogenomics knowledge discovery with the semantic web. *Brief Bioinform.* 2009;10(2):153–63.
- Coulet A, Garten Y, Dumontier M, Altman RB, Musen MA, Shah NH. Integration and publication of heterogeneous text-mined relationships on the Semantic Web. *J Biomed Semant.* 2011;2(S-2):S10.
- Samwald M, Giménez JM, Boyce RD, Freimuth RR, Adlassnig K-P, Dumontier M. Pharmacogenomic knowledge representation, reasoning and genome-based clinical decision support based on OWL 2 DL ontologies. *BMC Med Inform Dec Making.* 2015;15:12.
- Monnin P, Jonquet C, Legrand J, Napoli A, Coulet A. PGxO: A very lite ontology to reconcile pharmacogenomic knowledge units. In: *Methods, tools & platforms for Personalized Medicine in the Big Data Era. NETTAB 2017 Workshop Collection*. Palermo: PeerJ PrePrints; 2017. p. 1–4.
- Noy NF, McGuinness DL, et al. Ontology development 101: A guide to creating your first ontology. Stanford, CA: Stanford knowledge systems laboratory technical report KSL-01-05 and Stanford medical informatics technical report SMI-2001-0880; 2001.
- Dieng R, Corby O, Giboin A, Ribiere M. Methods and tools for corporate knowledge management. *Int J Hum Comput Stud.* 1999;51(3):567–98.
- Musen MA. The protégé project: a look back and a look forward. *AI Matters.* 2015;1(4):4–12.
- Baader F, Calvanese D, McGuinness DL, Nardi D, Patel-Schneider PF. *The Description Logic Handbook: Theory, Implementation, and Applications*. Cambridge: Cambridge University Press; 2003.
- Noy NF, Shah NH, Whetzel PL, Dai B, Dorf M, Griffith N, et al. BioPortal: ontologies and integrated data resources at the click of a mouse. *Nucleic Acids Res.* 2009;37(Web Server issue):W170–3.
- PgxO summary page on the ncbi bioportal. Available from: <https://bioportal.bioontology.org/ontologies/PGXO>. Accessed 30 July 2018.
- PractikPharma. PgxO page on github. Available from: <https://github.com/practikpharma/PGxO>. Accessed 30 July 2018.
- Matentzoglou N, Malone J, Mungall C, Stevens R. MIRO: guidelines for minimum information for the reporting of an ontology. *J Biomed Semant.* 2018;9(1):6:1–13. Available from: <https://doi.org/10.1186/s13326-017-0172-7>.
- Gangemi A. Ontology design patterns for semantic web content. In: *The Semantic Web - ISWC 2005, 4th International Semantic Web Conference, ISWC 2005, Galway, Ireland, November 6-10, 2005, Proceedings*. Springer; 2005. p. 262–76.
- Rindflesch TC, Kilicoglu H, Fiszman M, Roseblat G, Shin D. Semantic MEDLINE: an advanced information management application for biomedicine. *Inf Serv Use.* 2011;31(1-2):15–21.
- Tsuruoka Y, Miwa M, Hamamoto K, Tsujii J, Ananiadou S. Discovering and visualizing indirect associations between biomedical concepts. *Bioinformatics.* 2011;27(13):111–9.
- PractikPharma. Mappings from pgxo to so-pharm, po, phare and genomic cds. Available from: <https://github.com/practikpharma/PGxO/raw/master/mappings/mapp1.owl>. Accessed 30 July 2018.
- PractikPharma. Mappings from pgxo to mesh, nci and snomed ct. Available from: <https://github.com/practikpharma/PGxO/raw/master/mappings/mapp2.owl>. Accessed 30 July 2018.
- Bose R, Frew J. Lineage retrieval for scientific data processing: a survey. *ACM Comput Surv.* 2005;37:1–28.
- Lebo T, Sahoo S, McGuinness D, Belhajjame K, Cheney J, Corsar D, et al. PROV-O: The PROV Ontology. W3C Recommendation. 2013;30.
- Gangemi A. Ontology:dolce+dns ultralite - odp. Available from: http://ontologydesignpatterns.org/wiki/Ontology:DOLCE+DnS_Ultralite. Accessed 30 July 2018.
- Horrocks I, Patel-Schneider PF, Bechhofer S, Tsarkov D. Owl rules: A proposal and prototype implementation. *Web Semant.* 2005;3(1):23–40.
- Motik B, Sattler U, Studer R. Query answering for OWL-DL with rules. *J Web Sem.* 2005;3(1):41–60. Available from: <https://doi.org/10.1016/j.websem.2005.05.001>.
- Krötzsch M. OWL 2 profiles: An introduction to lightweight ontology languages. In: *Reasoning Web. Semantic Technologies for Advanced Query Answering - 8th International Summer School 2012, Vienna, Austria, September 3-8, 2012. Proceedings*. Springer; 2012. p. 112–83.
- Dalleau K, Marzougui Y, Da Silva S, Ringot P, Ndiaye NC, Coulet A. Learning from biomedical linked data to suggest valid pharmacogenes. *J Biomed Semant.* 2017;8(1):16.
- Landrum MJ, Lee JM, Riley GR, Jang W, Rubinstein WS, Church DM, et al. Clinvar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* 2013;42(D1):D980–5.
- Piñero J, Queralt-Rosinach N, Bravo À, Deu-Pons J, Bauer-Mehren A, Baron M, et al. Disgenet: a discovery platform for the dynamical exploration of human diseases and their genes. *Database.* 2015;2015.
- Wishart DS, Knox C, Guo AC, Cheng D, Shrivastava S, Tzur Dan, et al. Drugbank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic Acids Res.* 2007;36(suppl_1):D901–6.
- Kuhn M, Letunic I, Jensen LJ, Bork P. The sider database of drugs and side effects. *Nucleic Acids Res.* 2015;44(D1):D1075–9.
- Callahan A, Cruz-Toledo J, Ansell P, Dumontier M. Bio2rdf release 2: improved coverage, interoperability and provenance of life science linked data. In: *The Semantic Web: Semantics and Big Data, 10th International Conference, ESWC 2013, Montpellier, France, May 26-30, 2013. Proceedings*. Springer; 2013. p. 200–12.
- Wei C-H, Kao H-Y, Lu Z. Pubtator: a web-based text mining tool for assisting biocuration. *Nucleic Acids Res.* 2013;41(W1):W518–22.
- Stenetorp P, Pyysalo S, Topić G, Ohta T, Ananiadou S, Tsujii J. BRAT: a web-based tool for NLP-assisted text annotation. In: *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics. Association for Computational Linguistics; 2012. p. 102–7.*
- PractikPharma. Guidelines of our yet unpublished annotated corpus. Available from: https://github.com/practikpharma/PGxCorpus/raw/master/annotation_guidelines.pdf. Accessed 30 July 2018.

45. Collobert R, Weston J, Bottou L, Karlen M, Kavukcuoglu K, Kuksa P. Natural language processing (almost) from scratch. *J Mach Learn Res.* 2011;12(Aug):2493–537.
46. Quan C, Hua L, Sun X, Bai W. Multichannel convolutional neural network for biological relation extraction. *BioMed Res Int.* 2016;2016:.
47. Lebrecht R, Collobert R. Word embeddings through hellinger PCA. In: Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2014, April 26-30, 2014, Gothenburg, Sweden. The Association for Computer Linguistics; 2014. p. 482–90. Available from: <http://aclweb.org/anthology/E/E14/E14-1051.pdf>.
48. Neuraz A, Chouchana L, Malamut G, Le Beller C, Roche D, Beaune P, et al. Phenome-wide association studies on a quantitative trait: application to TPMT enzyme activity and thiopurine therapy in pharmacogenomics. *PLoS Comput Biol.* 2013;9(12):e1003405.
49. Ramirez AH, Shi Y, Schildcrout JS, Delaney JT, Xu H, Oetjens MT, et al. Predicting warfarin dosage in European-Americans and African-Americans using DNA samples linked to an electronic health record. *Pharmacogenomics.* 2012;13(4):407–18.
50. Mosley JD, Shaffer CM, Van Driest SL, Weeke PE, Wells QS, Karnes JH, et al. A genome-wide association study identifies variants in KCNIP4 associated with ACE inhibitor-induced cough. *Pharmacogenomics J.* 2016;16(3):231–7.
51. Van Driest SL, McGregor TL, Velez Edwards DR, Saville BR, Kitchner TE, Hebringer SJ, et al. Genome-Wide Association Study of Serum Creatinine Levels during Vancomycin Therapy. *PLoS ONE.* 2015;10(6):e0127791.
52. Wells QS, Veatch OJ, Fessel JP, Joon AY, Levinson RT, Mosley JD, et al. Genome-wide association and pathway analysis of left ventricular function after anthracycline exposure in adults. *Pharmacogenet Genomics.* 2017;27(7):247–54.
53. Kawai VK, Cunningham A, Vear SI, Van Driest SL, Oginni A, Xu H, et al. Genotype and risk of major bleeding during warfarin treatment. *Pharmacogenomics.* 2014;15(16):1973–83.
54. Feng Q, Wei WQ, Chung CP, Levinson RT, Bastarache L, Denny JC, et al. The effect of genetic variation in PCSK9 on the LDL-cholesterol response to statin therapy. *Pharmacogenomics J.* 2017;17(2):204–8.
55. Karnes JH, Cronin RM, Rollin J, Teumer A, Pouplard C, Shaffer CM, et al. A genome-wide association study of heparin-induced thrombocytopenia using an electronic medical record. *Thromb Haemost.* 2015;113(4):772–81.
56. Roden DM, Pulley JM, Basford MA, Bernard GR, Clayton EW, Balsler JR, et al. Development of a large-scale de-identified DNA biobank to enable personalized medicine. *Clin Pharmacol Ther.* 2008;84(3):362–9.
57. Denny JC, Van Driest SL, Wei WQ, Roden DM. The Influence of Big (Clinical) Data and Genomics on Precision Medicine and Drug Development. *Clin Pharmacol Ther.* 2018;103(3):409–18.
58. Smith B, Ceusters W, Klagges B, Köhler J, Kumar A, Lomax J, et al. Relations in biomedical ontologies. *Genome Biol.* 2005;6(5):R46.
59. Gottesman O, et al. The Electronic Medical Records and Genomics (eMERGE) Network: past, present, and future. *Genet Med.* 2013;15(10):761–71.
60. Jannot AS, Zapletal E, Avillach P, Mamzer MF, Burgun A, Degoulet P. The Georges Pompidou University Hospital Clinical Data Warehouse: A 8-years follow-up experience. *Int J Med Inform.* 2017;102:21–8.
61. Relling MV, Gardner EE, Sandborn WJ, Schmiegelow K, Pui CH, Yee SW, et al. Clinical pharmacogenetics implementation consortium guidelines for thiopurine methyltransferase genotype and thiopurine dosing: 2013 update. *Clin Pharmacol Ther.* 2013;1(4):324–5.
62. Bodenreider O. The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic Acids Res.* 2004;32(Database-Issue):267–270. Available from: <https://doi.org/10.1093/nar/gkh061>.
63. Humphreys BL, Lindberg DA, Schoolman HM, Barnett GO. The Unified Medical Language System: an informatics research collaboration. *J Am Med Inform Assoc.* 1998;5(1):1–11.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

