

# Common Foundations for SHACL, ShEx, and PG-Schema

Anonymous Author(s)

## ABSTRACT

Graphs have emerged as an important foundation for a variety of applications, including capturing and reasoning over factual knowledge, semantic data integration, social networks, and providing factual knowledge for machine learning algorithms. To formalise certain properties of the data and to ensure data quality, there is a need to describe the *schema* of such graphs. Because of the breadth of applications and availability of different data models, such as RDF and property graphs, both the Semantic Web and the database community have independently developed *graph schema languages*: SHACL, ShEx, and PG-Schema. Each language has its unique approach to defining constraints and validating graph data, leaving potential users in the dark about their commonalities and differences. In this paper, we provide formal, concise definitions of the *core components* of each of these schema languages. We employ a uniform framework to facilitate a comprehensive comparison between the languages and identify a common set of functionalities, shedding light on both overlapping and distinctive features of the three languages.

### ACM Reference Format:

Anonymous Author(s). 2024. Common Foundations for SHACL, ShEx, and PG-Schema. In *Proceedings of The Web Conference*. ACM, New York, NY, USA, 20 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Driven by the unprecedented growth of interconnected data, *graph-based data representations* have emerged as an expressive and versatile framework for modelling and analysing connections in data sets [46]. This rapid growth however, has led to a proliferation of diverse approaches, each with its own identity and perspective.

The two most prominent graph data models are *RDF* (Resource Description Framework) [14] and *Property Graphs* [9]. In RDF, data is modelled as a collection of triples, each consisting of a subject, predicate, and object. Such triples naturally represent either edges in a directed labelled graph (where the predicates represent relationships between nodes), or attributes-value pairs of nodes. That is, objects can both be entities or atomic (literal) values. In contrast, Property Graphs model data as nodes and edges, where both can have labels and records attached, allowing for a flexible representation of attributes directly on the entities and relationships.

Similarly to the different data models, we are also seeing different approaches towards *schema languages* for graph-structured data. Traditionally, in the Semantic Web community, schema and constraint languages have been *descriptive*, focusing on flexibility to accommodate varying structures. However, there has been a growing need for more *prescriptive* schemas that focus on *validation of data*. At the same time, in the Database community, schemas have traditionally been prescriptive but, since the rise of semi-structured data, the demand for descriptive schemas has been growing. Thus,

the philosophies of schemas in the two communities have been growing closer together.

For RDF, there are two main schema languages: SHACL (Shapes Constraint Language) [29], which is also a W3C recommendation, and ShEx (Shape Expressions) [43]. In the realm of Property Graphs, the current main approach is PG-Schema [2, 3]. The development processes of these languages have been quite different. For SHACL and ShEx, the formal semantics were only introduced after their initial implementations, echoing the evolution of programming languages. Indeed, an analysis of SHACL’s expressive power and associated decision problems appeared in the literature [6, 7, 34, 38–40] only after it was published as a W3C recommendation, leading up to a fully *recursive variant* of the language [1, 5, 12, 13, 39], whose semantics had been left undefined in the standard. A similar scenario occurred with ShEx, where formal analyses were only conducted in later phases [8, 48]. PG-Schema developed in the opposite direction. Here, a group of experts from industry and academia first defined the main ideas in a sequence of research papers [2, 3] and the implementation is expected to follow.

Since these three languages have been developed in different communities, in the course of different processes, it is no surprise that they are quite different. SHACL, ShEx, and PG-Schema use an array of diverse approaches for defining how their components work, ranging from *declarative* (formulas that *specify what to look for*) to *generative* (expressions that *generate the matching content*), and even combinations thereof. The bottom line is that we are left with three approaches to express a “schema for graph-structured data” that are very different at first glance.

As a group of authors coming from both the Semantic Web and Database communities, we believe that there is a *need for common understanding*. While the functionalities of schemas and constraints used in the two communities largely overlap, it is a daunting task to understand the essence of languages, such as SHACL, ShEx, and PG-Schema. In this paper, we therefore aim to shed light on the common aspects and the differences between these three languages.

Using a common framework, we provide crisp definitions of the main aspects of the languages. Because the languages operate on different data models, as a first step we introduce the *Common Graph Data Model*, a mathematical representation of data that *canonically embeds* both RDF graphs and Property Graphs (see Section 2, which also develops general common foundations). Precise abstractions of the languages themselves are presented in Sections 3 (SHACL), 4 (ShEx), and 5 (PG-Schema); in the Appendices we explain how and why we sometimes deviate from the original formalisms. Each of these sections contains examples to give readers an immediate intuition about what kinds of conditions each language can express. Then, in Section 6, we present the *Common Graph Schema Language (CoGSL)*, which consists of functionalities shared by them all.

Casting all three languages in a common framework has the immediate advantage that the reader can identify common functionalities *based on the syntax only*: on the one hand, we aim at giving the same semantics to schema language components that

117 syntactically look the same, and on the other hand, we can provide  
 118 examples of properties that distinguish the three languages using  
 119 simple syntactic constructs that are not part of the common core.  
 120 Aside from corner cases, properties expressed using constructs out-  
 121 side the common core are generally not expressible in all three  
 122 languages. By providing an understanding of fundamental differ-  
 123 ences and similarities between the three schema languages, we  
 124 hope to benefit both practitioners in choosing a schema language  
 125 fitting their needs, and researchers in studying the complexity and  
 126 expressiveness of schema languages.

## 128 2 FOUNDATIONS

129 In this section we present some material that we will need in the  
 130 subsequent sections, and define a data model that consists of com-  
 131 mon aspects of RDF and Property Graphs.

### 133 2.1 A common data model

134 When developing a common framework for SHACL, ShEx, and  
 135 PG-Schema, the first challenge is establishing a *common data model*,  
 136 since SHACL and ShEx work on RDF, whereas PG-Schema works  
 137 on Property Graphs. Rather than using a model that generalizes  
 138 both RDF and Property Graphs, we propose a simple model, called  
 139 *common graphs*, which we obtained by asking what, fundamentally,  
 140 are the *common aspects* of RDF and Property Graphs (Appendix A  
 141 gives more details on the distilling of common graphs).

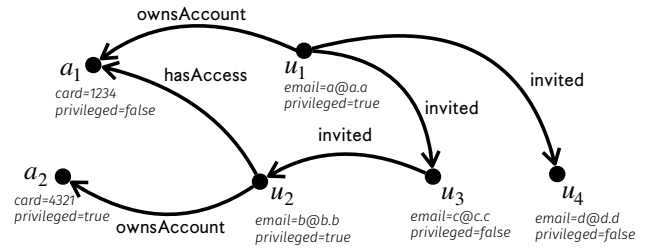
142 Let us assume disjoint countable sets of nodes  $\mathcal{N}$ , values  $\mathcal{V}$ ,  
 143 predicates  $\mathcal{P}$ , and keys  $\mathcal{K}$  (sometimes called properties).

144 **DEFINITION 1.** A common graph is a pair  $\mathcal{G} = (E, \rho)$  where

- 145 •  $E \subseteq_{\text{fin}} \mathcal{N} \times \mathcal{P} \times \mathcal{N}$  is its set of edges (which carry predicates), and
- 146 •  $\rho : \mathcal{N} \times \mathcal{K} \rightarrow \mathcal{V}$  is a finite-domain partial function mapping node-  
 147 key pairs to values.

148 The set of nodes of a common graph  $\mathcal{G}$ , written  $\text{Nodes}(\mathcal{G})$ , consists  
 149 of all elements of  $\mathcal{N}$  that occur in  $E$  or in the domain of  $\rho$ . Similarly,  
 150  $\text{Keys}(\mathcal{G})$  is the subset of  $\mathcal{K}$  that is used in  $\rho$ , and  $\text{Values}(\mathcal{G})$  is the  
 151 subset of  $\mathcal{V}$  that is used in  $\rho$  (that is, the range of  $\rho$ ).

152 **Example 1.** Consider Figure 1, containing a graph to store in-  
 153 formation about *users* who may have access to (possibly multiple)  
 154 *accounts* in, e.g., a media streaming service. In this example, we  
 155 have six nodes describing four persons ( $u_1, \dots, u_4$ ) and two accounts  
 156 ( $a_1, a_2$ ). As a common graph  $\mathcal{G} = (E, \rho)$ , the nodes are  $a_1, u_1$ , etc.  
 157 Examples of edges in  $E$  are  $(u_2, \text{hasAccess}, a_1)$  and  $(u_3, \text{invited}, u_2)$ .  
 158 Furthermore, we have  $\rho(u_2, \text{email}) = \text{d@d.d}$  and  $\rho(a_1, \text{card}) = 1234$ .  
 159 So,  $E$  captures the arrows in the figure (labeled with predicates)  
 160 and  $\rho$  captures the key/value information for each node. Notice  
 161 that a person may be the owner of an account, and may poten-  
 162 tially have access to other accounts. This is captured using the  
 163 predicates `ownsAccount` and `hasAccess`, respectively. In addition,  
 164 the system implements an invitation functionality, where users  
 165 may invite other people to join the platform. The previous invita-  
 166 tions are recorded using the predicate `invited`. Both accounts and  
 167 users may be privileged, which is stored via a Boolean value of the  
 168 key *privileged*. We note that the presence of the key *email* (resp.  
 169 of the key (credit) *card*) is associated with, and indeed identifies users  
 170 (resp. accounts).



175 **Figure 1: The media service common graph.**

176 It is easy to see that a common graph is a special case of a  
 177 property graph (see [2] for a formal definition of property graphs).  
 178 A common graph can also be seen as a set of triples, as in RDF. Let

$$179 \mathcal{E} = (\mathcal{N} \times \mathcal{P} \times \mathcal{N}) \cup (\mathcal{N} \times \mathcal{K} \times \mathcal{V}).$$

180 Then, a common graph can be seen as a finite set  $\mathcal{G} \subseteq \mathcal{E}$  such that  
 181 for each  $u \in \mathcal{N}$  and  $k \in \mathcal{K}$  there is at most one  $v \in \mathcal{V}$  such that  
 182  $(u, k, v) \in \mathcal{G}$ . Indeed, a common graph  $(E, \rho)$  corresponds to

$$183 E \cup \{(u, k, v) \mid \rho(u, k) = v\}.$$

184 When we write  $\rho(u, k) = v$  we assume that  $\rho$  is defined on  $(u, k)$ .

185 *Throughout the paper we see property graph  $\mathcal{G}$  simultaneously as  
 186 a pair  $(E, \rho)$  and as a set of triples from  $\mathcal{E}$ , switching between these  
 187 perspectives depending on what is most convenient at a given moment.*

### 190 2.2 Node contents and neighbourhoods

191 Let  $\mathcal{R}$  be the set of all *records*, i.e., finite-domain partial functions  
 192  $r : \mathcal{K} \rightarrow \mathcal{V}$ . We write records as sets of pairs  $\{(k_1, w_1), \dots, (k_n, w_n)\}$   
 193 where  $k_1, \dots, k_n$  are all different, meaning that  $k_i$  is mapped to  $w_i$ .

194 For a common graph  $\mathcal{G} = (E, \rho)$  and node  $v$  in  $\mathcal{G}$ , by a slight abuse  
 195 of notation we write  $\rho(v)$  for the record  $\{(k, w) \mid \rho(v, k) = w\}$  that  
 196 collects all key-value pairs associated with node  $v$  in  $\mathcal{G}$ . We call  
 197  $\rho(v)$  the *content* of node  $v$  in  $\mathcal{G}$ . This is how PG-Schema interprets  
 198 common graphs: it views key-value pairs in  $\rho(v)$  as *properties*  
 199 of the node  $v$ , rather than independent, navigable objects in the graph.

200 SHACL and ShEx, on the other hand, view common graphs as sets  
 201 of triples and make little distinction between keys and predicates.  
 202 The following notion—when applied to a node—uniformly captures  
 203 the local context of this node from that perspective: the content of  
 204 the node and all edges incident with the node.

205 **DEFINITION 2 (NEIGHBOURHOOD).** Given a common graph  $\mathcal{G}$   
 206 and a node or value  $v \in \mathcal{N} \cup \mathcal{V}$ , the neighbourhood of  $v$  in  $\mathcal{G}$  is  
 207  $\text{Neigh}_{\mathcal{G}}(v) = \{(u_1, p, u_2) \in \mathcal{G} \mid u_1 = v \text{ or } u_2 = v\}$ .

208 That is, when  $v \in \mathcal{N}$ , then  $\text{Neigh}_{\mathcal{G}}(v)$  is a star-shaped graph  
 209 where only the central node has non-empty content. When  $v \in$   
 210  $\mathcal{V}$ , then  $\text{Neigh}_{\mathcal{G}}(v)$  consists of all the nodes in  $\mathcal{G}$  that have some  
 211 key with value  $v$ , which is a common graph with no edges and a  
 212 restricted function  $\rho$ .

### 215 2.3 Value types

216 We assume an enumerable set of *value types*  $\mathcal{T}$ . The reader should  
 217 think of value types as integer, boolean, date, etc. Formally, for  
 218 each value type  $\mathfrak{v} \in \mathcal{T}$ , we assume that there is a set  $\llbracket \mathfrak{v} \rrbracket \subseteq \mathcal{V}$  of  
 219 all values of that type and that each value  $v \in \mathcal{V}$  belongs to some

type, i.e., there is at least one  $v \in \mathcal{T}$  such that  $v \in \llbracket v \rrbracket$ . Finally, we assume that there is a type  $\text{any} \in \mathcal{T}$  such that  $\llbracket \text{any} \rrbracket = \mathcal{V}$ .

## 2.4 Shapes and schemas

We formulate all three schema languages using *shapes*, which are unary formulas describing the graph's structure around a *focus* node or a value. Shapes will be expressed in different formalisms, specific to the schema language; for each of these formalisms we will define when a focus node or value  $v \in \mathcal{N} \cup \mathcal{V}$  *satisfies* shape  $\varphi$  in a common graph  $\mathcal{G}$ , written  $\mathcal{G}, v \models \varphi$ .

Inspired by ShEx *shape maps*, we abstract a schema  $\mathcal{S}$  as a set of pairs  $(sel, \varphi)$ , where  $\varphi$  is a shape and  $sel$  is a *selector*. A selector is also a shape, but usually a very simple one, typically checking the presence of an incident edge with a given predicate, or a property with a given key. A graph  $\mathcal{G}$  is *valid wrt.  $\mathcal{S}$* , written  $\mathcal{G} \models \mathcal{S}$ , if

$$\mathcal{G}, v \models sel \quad \text{implies} \quad \mathcal{G}, v \models \varphi$$

for all  $v \in \mathcal{N} \cup \mathcal{V}$  and  $(sel, \varphi) \in \mathcal{S}$ . That is, for each focus node or value satisfying the selector, the graph around it looks as specified by the shape. We call schemas  $\mathcal{S}$  and  $\mathcal{S}'$  *equivalent* if  $\mathcal{G} \models \mathcal{S}$  iff  $\mathcal{G} \models \mathcal{S}'$  for all  $\mathcal{G}$ . In what follows, we may use  $sel \Rightarrow \varphi$  to indicate a pair  $(sel, \varphi)$  from a schema  $\mathcal{S}$ .

*Example 2.* We next describe some constraints one may want to express in the domain of Example 1.

- (C1) We may want the values associated to certain keys to belong to concrete datatypes, like strings or Boolean values. In our example, we want to state that the value of the key *card* is always an integer.
- (C2) We may expect the existence of a value associated to a key, an outgoing edge, or even a complex path for a given source node. For our example, we require that all owners of an account have an email address defined.
- (C3) We may want to express database-like uniqueness constraints. For instance, we may wish to ensure that the email address of an account owner uniquely identifies them.
- (C4) We may want to ensure that all paths of a certain kind end in nodes with some desired properties. For example, if an account is privileged, then all users that have access to it should also be privileged.
- (C5) We may want to put an upper bound on the number of nodes reached from a given node by certain paths. For instance, every user may have access to at most 5 accounts.

## 3 SHACL ON COMMON GRAPHS

We first treat SHACL, because it is conceptually the simplest of the three languages. It is essentially a logic—some call it a *description logic in disguise* [6]. For each of the three languages, we need to perform some minor deviations in order to define it over common graphs. For SHACL, we discuss these in Appendix B. Our abstraction of SHACL on common graphs is inspired by [1, 6, 7, 13, 16].

**DEFINITION 3 (PATH EXPRESSION).** A path expression  $\pi$  is given by the following grammar:

$$\pi ::= \text{id} \mid p \mid k \mid \pi^- \mid \pi \cdot \pi \mid \pi \cup \pi \mid \pi^*$$

with  $p \in \mathcal{P}$ ,  $k \in \mathcal{K}$  and  $\text{id}$  the identity relation (or empty word).

**Table 1: Evaluation of a path expressions.**

$\pi$	$\llbracket \pi \rrbracket^{\mathcal{G}} \subseteq (\mathcal{N} \cup \mathcal{V}) \times (\mathcal{N} \times \mathcal{V})$ for $\mathcal{G} = (E, \rho)$
$\text{id}$	$\{(v, v) \mid v \in \mathcal{N} \cup \mathcal{V}\}$
$p$	$\{(v, u) \mid (v, p, u) \in E\}$
$k$	$\{(v, u) \mid \rho(v, k) = u\}$
$\pi^-$	$\{(v, u) \mid (u, v) \in \llbracket \pi \rrbracket^{\mathcal{G}}\}$
$\pi \cdot \pi'$	$\{(v, u) \mid \exists v' : (v, v') \in \llbracket \pi \rrbracket^{\mathcal{G}} \wedge (v', u) \in \llbracket \pi' \rrbracket^{\mathcal{G}}\}$
$\pi \cup \pi'$	$\llbracket \pi \rrbracket^{\mathcal{G}} \cup \llbracket \pi' \rrbracket^{\mathcal{G}}$
$\pi^*$	$\llbracket \text{id} \rrbracket^{\mathcal{G}} \cup \llbracket \pi \rrbracket^{\mathcal{G}} \cup \llbracket \pi \cdot \pi \rrbracket^{\mathcal{G}} \cup \dots$

**DEFINITION 4 (SHACL SHAPE).** A SHACL shape  $\varphi$  is given by the following grammar:

$$\varphi ::= \top \mid \text{is}(c) \mid \text{test}(v) \mid \text{closed}(Q) \mid \text{eq}(\pi, p) \mid \\ \text{disj}(\pi, p) \mid \neg\varphi \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \exists^{\geq n}\pi.\varphi \mid \exists^{\leq n}\pi.\varphi.$$

with  $c \in \mathcal{V}$ ,  $v \in \mathcal{T}$ ,  $Q \subseteq_{\text{fin}} \mathcal{P} \cup \mathcal{K}$ ,  $p \in \mathcal{P}$ , and  $n$  a natural number. We may use  $\exists\pi.\varphi$  as syntactic sugar for  $\exists^{\geq 1}\pi.\varphi$ .

**DEFINITION 5 (SHACL SELECTOR).** A SHACL selector  $sel$  is a SHACL shape of a restricted form, given by the following grammar:

$$sel ::= \exists p.\top \mid \exists k.\top \mid \exists p^-. \top \mid \exists k^-. \top \mid \text{is}(c).$$

with  $p \in \mathcal{P}$ ,  $k \in \mathcal{K}$ , and  $c \in \mathcal{V}$ .

Putting it together, a SHACL Schema  $\mathcal{S}$  is a finite set of pairs  $(sel, \varphi)$ , where  $sel$  is a SHACL selector and  $\varphi$  is a SHACL shape.

To define the semantics of SHACL schemas, we first define in Table 1 the semantics of a SHACL path expression  $\pi$  on a graph  $\mathcal{G}$  as a binary relation  $\llbracket \pi \rrbracket^{\mathcal{G}}$  over  $\mathcal{N} \cup \mathcal{V}$ . The semantics of SHACL shapes is defined in Table 2, which specifies when a node or value  $v$  *satisfies* a SHACL shape  $\varphi$  w.r.t. a  $\mathcal{G}$ , written  $\mathcal{G}, v \models \varphi$ . Note that both  $\llbracket \pi \rrbracket^{\mathcal{G}}$  and  $\{v \in \mathcal{N} \cup \mathcal{V} \mid \mathcal{G}, v \models \varphi\}$  may be infinite: for example,  $\llbracket \text{id} \rrbracket^{\mathcal{G}}$  is the identity relation over the infinite set  $\mathcal{N} \cup \mathcal{V}$ .

The semantics of SHACL schemas then follows Section 2.4. Importantly, SHACL selectors always select a finite subset of  $\mathcal{N} \cup \mathcal{V}$ : the selected nodes or values come either from the selector itself, in the case of  $\text{is}(c)$ , or from  $\mathcal{G}$ , in the remaining four cases. For example,  $\exists p.\top$  selects those nodes of  $\mathcal{G}$  that have an outgoing  $p$ -edge in  $\mathcal{G}$ —it is grounded to  $\mathcal{G}$  in the second line of Table 1. In consequence, each pair  $(sel, \varphi)$  in a SHACL schema tests the inclusion of a finite set of nodes or values in a possibly infinite set.

*Example 3.* For better readability we write  $\exists\pi$  instead of  $\exists^{\geq 1}\pi.\top$  (that is, we omit  $\top$ ) and  $\forall\pi.\varphi$  instead of  $\exists^{\leq 0}\pi.\neg\varphi$ . Let us see how the constraints from Example 2 can be handled in SHACL. For (C1), we assume the value type  $\text{int}$  with the obvious meaning. The following SHACL constraints express the constraints (C1–C5):

$$\exists \text{card}^- \Rightarrow \text{test}(\text{int}) \quad (\text{C1})$$

$$\exists \text{ownsAccount} \Rightarrow \exists \text{email} \quad (\text{C2})$$

$$\exists \text{email}^- \Rightarrow \exists^{\leq 1} \text{email}^- \quad (\text{C3})$$

$$\exists \text{card} \Rightarrow (\exists \text{privileged}.\neg \text{is}(\text{true})) \vee \\ \forall \text{hasAccess}^- . (\exists \text{privileged}.\text{is}(\text{true})) \quad (\text{C4})$$

**Table 2: Semantics of a SHACL shape  $\varphi$ .**

$\varphi$	$\mathcal{G}, v \models \varphi$ if:
$\top$	trivially satisfied
$\neg\varphi$	not $\mathcal{G}, v \models \varphi$
$\varphi \wedge \varphi'$	$\mathcal{G}, v \models \varphi$ and $\mathcal{G}, v \models \varphi'$
$\varphi \vee \varphi'$	$\mathcal{G}, v \models \varphi$ or $\mathcal{G}, v \models \varphi'$
$\text{is}(c)$	$v = c$
$\text{test}(\mathbb{v})$	$v \in \llbracket \mathbb{v} \rrbracket$
$\text{closed}(Q)$	$\forall p \in (\mathcal{P} \cup \mathcal{K}) \setminus Q : \text{not } \mathcal{G}, v \models \exists^{\geq 1} p. \top$
$\text{eq}(\pi, p)$	$\{u \mid (v, u) \in \llbracket \pi \rrbracket^{\mathcal{G}}\} = \{u \mid (v, u) \in \llbracket p \rrbracket^{\mathcal{G}}\}$
$\text{disj}(\pi, p)$	$\{u \mid (v, u) \in \llbracket \pi \rrbracket^{\mathcal{G}}\} \cap \{u \mid (v, u) \in \llbracket p \rrbracket^{\mathcal{G}}\} = \emptyset$
$\exists^{\geq n} \pi. \varphi$	$\#\{u \mid (v, u) \in \llbracket \pi \rrbracket^{\mathcal{G}} \wedge \mathcal{G}, u \models \varphi\} \geq n$
$\exists^{\leq n} \pi. \varphi$	$\#\{u \mid (v, u) \in \llbracket \pi \rrbracket^{\mathcal{G}} \wedge \mathcal{G}, u \models \varphi\} \leq n$

$$\exists \text{email} \Rightarrow \exists^{\leq 5} \text{hasAccess}. \quad (\text{C5})$$

Concerning constraint (C3), notice that by using inverse *email* edges, the constraint indeed states that the email addresses uniquely identify users.

The constructs  $\text{eq}(\pi, p)$  and  $\text{disj}(\pi, p)$  are unique to SHACL. Let us see them in use.

*Example 4.* Using  $\text{eq}(\pi, p)$ , we can say, for instance, that an owner of an account also has access to it:

$$\exists \text{ownsAccount} \Rightarrow \text{eq}(\text{hasAccess} \cup \text{ownsAccount}, \text{hasAccess}).$$

Note how we use  $\text{eq}$  and  $\cup$  to express that the existence of one path (*ownsAccount*) implies the existence of another path (*hasAccess*) with the same endpoints.

A key feature in SHACL that is not available in ShEx is the ability to use regular expressions to talk about complex paths. This provides a limited form of recursive navigation in the graph, even though the standard SHACL does not support recursive constraints (in contrast to standard ShEx). See below for an example.

*Example 5.* Suppose that in Figure 1, we wanted to express that a privileged user may only invite other privileged users, who in turn can also only invite other privileged users. One way to express this in SHACL is as follows:

$$\begin{aligned} \exists \text{privileged} \Rightarrow \exists \text{privileged}. \text{is}(\text{false}) \vee \\ \forall (\text{invited}^* \cdot \text{privileged}). \text{is}(\text{true}). \end{aligned}$$

## 4 SHEX ON COMMON GRAPHS

While SHACL is conceptually the simplest of the three language, ShEx lies at the opposite end of the spectrum. It is an intricate, mutually recursive combination of a simple logic for shapes and a powerful formalism (triple expressions) for generating the allowed neighbourhoods. In this work we consider non-recursive ShEx, where shapes and triple expressions can be nested multiple times, but cannot be actually recursive. We choose non-recursive ShEx because it is significantly easier to understand (so it aligns with our overall understandability goal). The abstraction of ShEx over

common graphs is based on the treatment of ShEx on RDF triples by Boneva et al. [8]. The correspondence to standard ShEx is discussed in Appendix C.

**DEFINITION 6 (SHAPES AND TRIPLE EXPRESSIONS).** *ShEx shapes  $\varphi$ , triple expressions  $e$ , and closed triple expressions  $f$  are defined by the following grammar*

$$\begin{aligned} \varphi &::= \text{is}(c) \mid \text{test}(\mathbb{v}) \mid \{e\} \mid \varphi \wedge \varphi' \mid \varphi \vee \varphi' \mid \neg\varphi. \\ e &::= f; \text{op}_- \mid f; \text{op}_\pm \\ \text{op}_- &::= (\neg P^-)^* \\ \text{op}_\pm &::= (\neg P^-)^*; (\neg Q)^* \\ f &::= \varepsilon \mid p. \varphi \mid p^- . \varphi \mid f; f \mid f \mid f \mid f^*. \end{aligned}$$

where  $c \in \mathcal{V}$ ,  $\mathbb{v} \in \mathcal{T}$ ,  $p \in \mathcal{P} \cup \mathcal{K}$ , and  $P, Q \subseteq_{\text{fin}} \mathcal{P} \cup \mathcal{K}$ .

The notion of satisfaction for ShEx shapes and the semantics of triple expressions are defined by mutual recursion in Table 3 and Table 4. Triple expressions are used to specify neighbourhoods of nodes and values. They require to consider incoming and outgoing edges separately. For this purpose we decorate incoming edges with  $^-$ . Formally, we introduce a fresh predicate  $p^-$  for each  $p \in \mathcal{P}$  and a fresh key  $k^-$  for each  $k \in \mathcal{K}$ . We let  $\mathcal{P}^- = \{p^- \mid p \in \mathcal{P}\}$ ,  $\mathcal{K}^- = \{k^- \mid k \in \mathcal{K}\}$ ,  $\mathcal{E}^- = \mathcal{N} \times \mathcal{P}^- \times \mathcal{N} \cup \mathcal{V} \times \mathcal{K}^- \times \mathcal{N}$ , and define  $\text{Neigh}_{\mathcal{G}}^{\pm}(v) \subseteq \mathcal{E} \cup \mathcal{E}^-$  as

$$\{(v, p, v') \mid (v, p, v') \in \mathcal{G}\} \cup \{(v, p^-, v') \mid (v', p, v) \in \mathcal{G}\}.$$

Compared to  $\text{Neigh}_{\mathcal{G}}(v)$ , apart from flipping the incoming edges and marking them with  $^-$ , we also represent each loop  $(v, p, v)$  twice: once as an outgoing edge  $(v, p, v)$  and once as an incoming edge  $(v, p^-, v)$ . Notice that in Table 4,  $\neg P^-$  and its inverse counterpart  $\neg P^-$  are treated as any other triple expressions, even though in the grammar they are allowed only at the top level.

Closed triple expressions  $f$  define neighbourhoods that use only a finite number of predicates and keys; such neighbourhoods are also called *closed* in ShEx terminology. General triple expressions  $e$  *open* the neighbourhood either only w.r.t. incoming triples ( $\text{op}_-$ ) by allowing any incoming triples whose predicate or key is not in a set  $P$ , or w.r.t. both incoming and outgoing triples ( $\text{op}_\pm$ ) by additionally allowing outgoing triples whose predicate or key is not in a set  $Q$ . Let  $\top = \varepsilon; (-\emptyset^-)^*; (-\emptyset)^*$ . Then  $\top$  describes all possible neighbourhoods and  $\{\top\}$  is satisfied in every node and in every value of every graph.

**DEFINITION 7 (SHEX SELECTORS).** *A ShEx selector is a ShEx shape of a restricted form, defined by the grammar*

$$\begin{aligned} \text{sel} &::= \text{is}(c) \mid \{p. \text{is}(c); \top\} \mid \{p^- . \text{is}(c); \top\} \mid \\ &\quad \{p. \{\top\}; \top\} \mid \{p^- . \{\top\}; \top\}. \end{aligned}$$

where  $p \in \mathcal{P} \cup \mathcal{K}$  and  $c \in \mathcal{V}$ .

Following Section 2.4, a *ShEx schema*  $\mathcal{S}$  is a set of pairs of the form  $(\text{sel}, \varphi)$  where  $\varphi$  is a ShEx shape and  $\text{sel}$  is a ShEx selector.

We will be using these syntactic short-hands:

- $e^n$  for some positive integer  $n$  denotes the triple expression  $e; \dots; e$  where  $e$  is used  $n$  times,
- $e^{\leq n}$  as a short-hand for  $\varepsilon \mid e^1 \mid \dots \mid e^n$  and  $e^{\geq n}$  for  $e^n; e^*$ ,
- for a closed triple expression  $f$ , we let

$$\{f\}^\circ = \{f; (\neg P^-)^*; (\neg P)^*\}$$

**Table 3: Satisfaction of ShEx shapes.**

$\varphi$	$\mathcal{G}, v \models \varphi$ for $v \in \mathcal{N} \cup \mathcal{V}$
$test(\mathbb{v})$	$v \in \llbracket \mathbb{v} \rrbracket$
$is(c)$	$v = c$
$\{e\}$	$Neigh_{\mathcal{G}}^{\pm}(v) \in \llbracket e \rrbracket_v^{\mathcal{G}}$
$\varphi_1 \wedge \varphi_2$	$\mathcal{G}, v \models \varphi_1$ and $\mathcal{G}, v \models \varphi_2$
$\varphi_1 \vee \varphi_2$	$\mathcal{G}, v \models \varphi_1$ or $\mathcal{G}, v \models \varphi_2$
$\neg\varphi$	not $\mathcal{G}, v \models \varphi$

**Table 4: Semantics of triple expressions.**

$e$	$\llbracket e \rrbracket_v^{\mathcal{G}} \subseteq \mathcal{E} \cup \mathcal{E}^-$
$\varepsilon$	$\{\emptyset\}$
$p.\varphi$	$\{(v, p, v') \mid \mathcal{G}, v' \models \varphi\}$
$p^-. \varphi$	$\{(v, p^-, v') \mid \mathcal{G}, v' \models \varphi\}$
$\neg P$	$\{(v, p, v') \mid \mathcal{G}, v' \not\models P\}$
$\neg P^-$	$\{(v, p^-, v') \mid \mathcal{G}, v' \not\models P\}$
$e_1 ; e_2$	$\{T_1 \cup T_2 \mid T_1 \in \llbracket e_1 \rrbracket_v^{\mathcal{G}}, T_2 \in \llbracket e_2 \rrbracket_v^{\mathcal{G}}, T_1 \cap T_2 = \emptyset\}$
$e_1   e_2$	$\llbracket e_1 \rrbracket_v^{\mathcal{G}} \cup \llbracket e_2 \rrbracket_v^{\mathcal{G}}$
$e^*$	$\{\emptyset\} \cup \bigcup_{n=1}^{\infty} \left\{ T_1 \cup \dots \cup T_n \mid \begin{array}{l} T_1, \dots, T_n \in \llbracket e \rrbracket_v^{\mathcal{G}} \text{ and} \\ T_i \cap T_j = \emptyset \text{ for all } i \neq j \end{array} \right\}$

where  $P$  is the set of predicates and keys that appear *directly* in  $f$  without considering those that appear in  $\varphi$  for a sub-expression of the form  $p.\varphi$ . For instance, if  $f = p.\{p'.is(c)\}$ , then  $P = \{p\}$ . Also,  $P^-$  is the set of predicates and keys that appear inversed in  $f$ .

*Example 6.* Let us now see how the concrete constraints from Example 2 can be handled in ShEx.

$$\{card^-. \{ \top ; \top \} \Rightarrow test(\text{imt}) \quad (C1)$$

$$\{ownsAccount. \{ \top ; \top \} \Rightarrow \{email. \{ \top \} \}^\circ \quad (C2)$$

$$\{email^-. \{ \top ; \top \} \Rightarrow \{ \{email^-. \{ \top \} \}^{\leq 1} \}^\circ \quad (C3)$$

$$\{card. \{ \top ; \top \} \Rightarrow \{ \{privileged.^- . is(true) \}^\circ \}^* \vee \{ \{ (hasAccess^-. \{privileged.is(true) \}^\circ)^* \}^\circ \quad (C4)$$

$$\{email. \{ \top ; \top \} \Rightarrow \{ \{ (hasAccess. \{ \top \} )^{\leq 5} \}^\circ \quad (C5)$$

We next show a more complex example, which illustrates the power of ShEx that is not readily available in SHACL or PG-Schema.

*Example 7.* Suppose that we want to express the following constraint on each user who owns an account: the number of accounts to which the user has access is greater or equal to the number of accounts that the user owns. We can do this in ShEx as follows:

$$\{ownsAccount. \{ \top ; \top \} \Rightarrow \{ (hasAccess. \{ \top \} )^* ; (ownsAccount. \{ \top ; \top \} ; hasAccess. \{ \top \} )^* \}^\circ$$

Finally, let us see why ShEx and SHACL count differently.

*Example 8 (ShEx counts edges).* The following SHACL schema expresses that from every node with an outgoing hasAccess-edge,

**Table 5: Semantics of content types.**

$\mathbb{c}$	$\llbracket \mathbb{c} \rrbracket \subseteq \mathcal{R}$
$\llbracket \top \rrbracket$	$\mathcal{R}$
$\llbracket \{ \} \rrbracket$	$\{\mathbf{r}_\emptyset\}$
$\llbracket \{k : \mathbb{v}\} \rrbracket$	$\{ \{(k, w)\} \mid w \in \llbracket \mathbb{v} \rrbracket \}$
$\llbracket \mathbb{c}_1 \& \mathbb{c}_2 \rrbracket$	$\{ (r_1 \cup r_2) \in \mathcal{R} \mid r_1 \in \llbracket \mathbb{c}_1 \rrbracket \wedge r_2 \in \llbracket \mathbb{c}_2 \rrbracket \}$
$\llbracket \mathbb{c}_1   \mathbb{c}_2 \rrbracket$	$\llbracket \mathbb{c}_1 \rrbracket \cup \llbracket \mathbb{c}_2 \rrbracket$

there should be exactly two nodes accessible via a hasAccess-edge or an ownsAccount-edge:

$$\exists hasAccess \Rightarrow \exists^{=2} (hasAccess \cup ownsAccount). \top ;$$

here  $\exists^{=n} \pi.\varphi$  is a shorthand for  $\exists^{\leq n} \pi.\varphi \wedge \exists^{\geq n} \pi.\varphi$ . For instance, the graph below on the left is valid, whereas the one on the right is not.



The same constraint cannot be expressed in ShEx because ShEx cannot distinguish these two graphs. Indeed, ShEx triple expressions count triples adjacent to a node, whereas SHACL and PG-Schema count nodes on the opposite end of such triples.

## 5 PG-SCHEMA ON COMMON GRAPHS

PG-Schema is a non-recursive combination of a logic and two generative formalisms. It uses path expressions to specify paths (as in SHACL), and *content types* to specify node contents. Both path expressions and content types are then used in formulas defining shapes. Content types in PG-Schema play a role similar to triple expressions in ShEx, but they are only used for properties. Because all properties of a node must have different keys, they are much simpler than triple expressions (in fact, they can be translated into a fragment of SHACL). Unlike for SHACL and ShEx, the abstraction of PG-Schema on common graphs departs significantly from the original design. Original PG-Schema uses queries written in an external query language, which is left unspecified aside from some basic assumptions about the expressive power. Here we use a specific query language (PG-path expressions). Importantly, up to the choice of the query language, the abstraction we present here faithfully captures the expressive power of the original PG-Schema. A detailed comparison can be found in Appendix D.

**DEFINITION 8 (CONTENT TYPE).** A content type is an expression  $\mathbb{c}$  of the form defined by the grammar

$$\mathbb{c} ::= \top \mid \{ \} \mid \{k : \mathbb{v}\} \mid \mathbb{c} \& \mathbb{c} \mid \mathbb{c} | \mathbb{c} .$$

where  $k \in \mathcal{K}$  and  $\mathbb{v} \in \mathcal{T}$ .

Recall that  $\mathcal{R}$  is the set of all records (finite-domain partial functions  $r : \mathcal{K} \rightarrow \mathcal{V}$ ). We write  $\mathbf{r}_\emptyset$  for the empty record. For records  $r_1$  and  $r_2$ , we let  $r_1 \cup r_2$  be the function that behaves as  $r_1$  on  $\text{dom}(r_1)$  and as  $r_2$  on  $\text{dom}(r_2)$ . We require that  $r_1(k) = r_2(k)$  for every  $k \in \text{dom}(r_1) \cap \text{dom}(r_2)$ . The semantics of content types is defined in Table 5. Note that  $\llbracket \mathbb{c} \rrbracket$  is independent from  $\mathcal{G}$  and can be infinite.

**Table 6: Semantics of PG-path expressions.**

$\pi$	$\llbracket \pi \rrbracket^{\mathcal{G}} \subseteq (\mathcal{N} \cup \mathcal{V}) \times (\mathcal{N} \cup \mathcal{V})$ for $\mathcal{G} = (E, \rho)$
$\{k : c\}$	$\{(u, u) \mid u \in \text{Nodes}(\mathcal{G}) \wedge (k, c) \in \rho(u)\}$
$\neg\{k : c\}$	$\{(u, u) \mid u \in \text{Nodes}(\mathcal{G}) \wedge (k, c) \notin \rho(u)\}$
$\mathbb{C}$	$\{(u, u) \mid u \in \text{Nodes}(\mathcal{G}) \wedge \rho(u) \in \llbracket \mathbb{C} \rrbracket\}$
$\neg\mathbb{C}$	$\{(u, u) \mid u \in \text{Nodes}(\mathcal{G}) \wedge \rho(u) \notin \llbracket \mathbb{C} \rrbracket\}$
$k$	$\{(u, w) \mid \rho(u, k) = w\}$
$p$	$\{(u, v) \mid (u, p, v) \in E\}$
$\neg P$	$\{(u, v) \mid \exists p : (u, p, v) \in E \wedge p \notin P\}$
$\pi^-$	$\{(u, v) \mid (v, u) \in \llbracket \pi \rrbracket^{\mathcal{G}}\}$
$\pi \cdot \pi'$	$\{(u, v) \mid \exists w : (u, w) \in \llbracket \pi \rrbracket^{\mathcal{G}} \wedge (w, v) \in \llbracket \pi' \rrbracket^{\mathcal{G}}\}$
$\pi \cup \pi'$	$\llbracket \pi \rrbracket^{\mathcal{G}} \cup \llbracket \pi' \rrbracket^{\mathcal{G}}$
$\pi^*$	$\{(u, u) \mid u \in \text{Nodes}(\mathcal{G})\} \cup \llbracket \pi \rrbracket^{\mathcal{G}} \cup \llbracket \pi \cdot \pi \rrbracket^{\mathcal{G}} \cup \dots$

*Example 9.* We assume integers and strings are represented via `int`, `str`  $\in \mathcal{T}$ . Suppose we want to create a content type for nodes that have a string value for the *email* key and *optionally* have an integer value for the *card* key. No other key-value pairs are allowed. We should then use  $\{email : str\} \& (\{card : int\} \mid \{\})$ .

**DEFINITION 9 (PG-PATH EXPRESSIONS).** A PG-path expression is an expression  $\pi$  of the form defined by the grammar

$$\begin{aligned} \pi &::= \bar{\pi} \mid \bar{\pi} \cdot k \mid k^- \cdot \bar{\pi} \mid k^- \cdot \bar{\pi} \cdot k' \\ \bar{\pi} &::= \{k : c\} \mid \neg\{k : c\} \mid \mathbb{C} \mid \neg\mathbb{C} \mid p \mid \neg P \mid \bar{\pi}^- \mid \bar{\pi} \cdot \bar{\pi} \mid \bar{\pi} \cup \bar{\pi} \mid \bar{\pi}^* \end{aligned}$$

where  $k, k' \in \mathcal{K}$ ,  $c \in \mathcal{V}$ ,  $\mathbb{C}$  is a content type,  $p \in \mathcal{P}$ , and  $P \subseteq_{fin} \mathcal{P}$ . We use  $k, k^-$ , and  $k^- \cdot k'$  as short-hands for PG-path expressions  $\top \cdot k, k^- \cdot \top$ , and  $k^- \cdot \top \cdot k'$ , respectively.

Unlike in SHACL, PG-path expressions cannot navigate freely through values. In the property graph world, this would correspond to a join, which is a costly operation. Indeed, existing query languages for property graphs do not allow joins under  $*$ . However, PG-path expressions can start in a value and finish in a value. This leads to *node-to-node*, *node-to-value*, *value-to-node*, and *value-to-value* PG-path expressions, reflected in the four cases in the first rule of the grammar.

The semantics of PG-path expression  $\pi$  for graph  $\mathcal{G}$  is a binary relation over  $\text{Nodes}(\mathcal{G}) \cup \text{Values}(\mathcal{G})$ , defined in Table 6. In the table,  $k$  is treated as any other subexpressions, eventhough it can only be used at the end of a PG-path expression, or in the beginning as  $k^-$ . Notice that  $\neg\mathbb{C}$  matches nodes whose content is not of type  $\mathbb{C}$ ,  $\neg P$  matches edges with a label that is not in  $P$  (in particular,  $\neg\emptyset$  matches all edges). Also,  $\llbracket \pi \rrbracket^{\mathcal{G}}$  is always a subset of  $\mathcal{N} \times \mathcal{N}$ ,  $\mathcal{N} \times \mathcal{V}$ ,  $\mathcal{V} \times \mathcal{N}$ , or  $\mathcal{V} \times \mathcal{V}$ , corresponding to the four kinds of PG-path expressions discussed above.

**DEFINITION 10 (PG-SHAPES).** A PG-Shape is an expression  $\varphi$  defined by the following grammar:

$$\varphi ::= \exists^{\leq n} \pi \mid \exists^{\geq n} \pi \mid \varphi \wedge \varphi$$

where  $\pi$  is a PG-path expression. We use  $\exists$  and  $\exists^{\geq 1}$  as short-hands for  $\exists^{\geq 1}$  and  $\exists^{\leq 0}$ .

**Table 7: Satisfaction of PG-shapes**

$\varphi$	$\mathcal{G}, v \models \varphi$ for $v \in \mathcal{N} \cup \mathcal{V}$
$\exists^{\leq n} \pi$	$\#\{v' \mid (v, v') \in \llbracket \pi \rrbracket^{\mathcal{G}}\} \leq n$
$\exists^{\geq n} \pi$	$\#\{v' \mid (v, v') \in \llbracket \pi \rrbracket^{\mathcal{G}}\} \geq n$
$\varphi_1 \wedge \varphi_2$	$\mathcal{G}, v \models \varphi_1$ and $\mathcal{G}, v \models \varphi_2$

The semantics of PG-shapes is defined in Table 7. We say  $v \in \mathcal{N} \cup \mathcal{V}$  satisfies a PG-shape  $\varphi$  in a graph  $\mathcal{G}$  if  $\mathcal{G}, v \models \varphi$ . Every PG-shape is satisfied by nodes only or by values only.

**DEFINITION 11 (PG-SELECTORS).** A PG-selector is a PG-shape of the form  $\exists \pi$ .

A PG-Schema  $\mathcal{S}$  is a finite set of pairs  $(sel, \varphi)$  where  $sel$  is a PG-selector and  $\varphi$  is a PG-shape. The semantics of PG-Schemas is defined just like in Section 2.4.

*Example 10.* The constraints (C1-C5) from Example 2 can be handled in PG-Schema as follows:

$$\exists card^- \Rightarrow \exists(\{card : int\} \& \top) \quad (C1)$$

$$\exists ownsAccount \Rightarrow \exists email \quad (C2)$$

$$\exists email^- \Rightarrow \exists^{\leq 1} email^- \quad (C3)$$

$$\begin{aligned} \exists(\{card : any\} \& \top) \cdot \{\text{privileged} : true\} \Rightarrow \\ \nexists hasAccess^- \cdot \neg\{\text{privileged} : true\} \end{aligned} \quad (C4)$$

$$\exists email \Rightarrow \exists^{\leq 5} hasAccess \quad (C5)$$

A characteristic feature of PG-Schema, revealing its database provenience, is that it can close the whole graph by imposing restrictions on all nodes.

*Example 11.* Given a common graph such as the one in Figure 1, we might want to express that each node has a key *privileged* with a boolean value and either a key *card* with an integer value or a key *email* with a string value, and no other keys are allowed. In PG-Schema this can be expressed as follows:

$$\exists \top \Rightarrow \exists(\{\text{privileged} : bool\} \& (\{card : int\} \mid \{email : str\}))$$

We can also forbid any predicates except those mentioned in the running example:

$$\exists \top \Rightarrow \nexists \neg\{ownsAccount, hasAccess, invited\}$$

## 6 COMMON GRAPH SCHEMA LANGUAGE

We now present the Common Graph Schema Language (CoGSL), which combines the core functionalities shared by SHACL, ShEx, and PG-Schema (over common graphs).

Let us begin by examining the restrictions that need to be imposed. We shall refer to shapes and selectors used in CoGSL as *common shapes* and *common selectors*.

Common shapes cannot be closed under disjunction and negation, because PG-Schema shapes are purely conjunctive. For the same reason common shapes cannot be nested.

Kleene star  $*$  cannot be allowed in path expressions because we consider ShEx without recursion. By switching to ShEx with recursion, we would be able to support arbitrary SHACL path expressions

in shapes of the form  $\exists\pi$ , but not arbitrary PG-path expressions as these are too expressive for SHACL.

Supporting path expressions traversing more than one edge under counting quantifiers is impossible as this would not be expressible in ShEx. Supporting disjunctions of labels of the form  $p_1 \cup p_2$  is also impossible, due to a mismatch in the approach to counting: while SHACL and PG-Schema count nodes and values, ShEx counts triples, as illustrated in Example 8.

Closed content types and  $\neg P$  cannot be used freely, because neither SHACL nor ShEx are capable of closing only properties or only predicate edges: both must be closed at the same time.

Finally, selectors are restricted because SHACL and ShEx do not support  $\top$  as a selector; that is, one cannot say that each node (or value) in the graph satisfies a given shape. This means that SHACL and ShEx schemas always allow a disconnected part of the graph that uses only predicates and keys not mentioned in the schema, whereas PG-Schema can disallow it (see Example 11).

Putting these restrictions together we obtain the Common Graph Schema Language. We define it below as a fragment of PG-Schema.

**DEFINITION 12 (COMMON SHAPE).** A common shape  $\varphi$  is an expression given by the grammar

$$\begin{aligned} \varphi &::= \exists\pi \mid \exists^{\leq n}\pi_1 \mid \exists^{\geq n}\pi_1 \mid \exists\mathbb{C} \wedge \nexists\neg P \mid \varphi \wedge \varphi . \\ \mathbb{C} &::= \{\} \mid \{k : \mathbb{V}\} \mid \mathbb{C} \& \mathbb{C} \mid \mathbb{C} \mid \mathbb{C} . \\ \pi_0 &::= \{k : c\} \mid \neg\{k : c\} \mid \mathbb{C} \& \top \mid \neg(\mathbb{C} \& \top) \mid \pi_0 \cdot \pi_0 . \\ \pi_1 &::= \pi_0 \cdot p \cdot \pi_0 \mid \pi_0 \cdot p^- \cdot \pi_0 \mid \pi_0 \cdot k \mid k^- \cdot \pi_0 . \\ \bar{\pi} &::= \pi_0 \mid p \mid \bar{\pi}^- \mid \bar{\pi} \cdot \bar{\pi} \mid \bar{\pi} \cup \bar{\pi} . \\ \pi &::= \bar{\pi} \mid \bar{\pi} \cdot k \mid k^- \cdot \bar{\pi} \mid k^- \cdot \bar{\pi} \cdot k' . \end{aligned}$$

where  $n \in \mathbb{N}$ ,  $P \subseteq_{fin} \mathcal{P}$ ,  $k, k' \in \mathcal{K}$ ,  $c \in \mathcal{V}$ , and  $p \in \mathcal{P}$ .

That is,  $\mathbb{C}$  is a content type that does not use  $\top$  (a *closed* content type),  $\pi_0$  is a PG-path expression that always stays in the same node (a *filter*),  $\pi_1$  is a PG-path expression that traverses a single edge or property (forward or backwards), and  $\pi$  is a PG-path expression that uses neither  $*$  nor  $\neg P$ . Moreover,  $\pi_0$ ,  $\pi_1$ , and  $\pi$  can only use *open* content types; that is, content types of the form  $\mathbb{C} \& \top$ . The use of  $\neg P$  is limited to closing the neighbourhood of a node (this is the only way PG-Schema can do it).

**DEFINITION 13 (COMMON SELECTOR).** A common selector is a common shape of one of the following forms

$$\exists k, \exists p \cdot \pi, \exists p^- \cdot \pi, \exists \{k : c\} \cdot \pi, \exists (\{k : \mathbb{V}\} \& \top) \cdot \pi, \exists k^- \cdot \pi,$$

where  $k \in \mathcal{K}$ ,  $p \in \mathcal{P}$ ,  $c \in \mathcal{V}$ ,  $\mathbb{V} \in \mathcal{T}$  and  $\pi = \bar{\pi}$  or  $\pi = \bar{\pi} \cdot k'$  for some PG-path expression  $\bar{\pi}$  generated by the grammar in Definition 12 and some  $k' \in \mathcal{K}$ .

That is, a common selector is a common shape of the form  $\exists\pi$  such that the PG-path expression  $\pi$  requires the focus node or value to occur in a triple with a specified predicate or key.

A *common schema* is a finite set of pairs  $(sel, \varphi)$  where  $sel$  is a common selector and  $\varphi$  is a common shape. The semantics is inherited from PG-Schema.

We note that we showed that the constraints (C1)-(C5) from our running example can be expressed in all three formalisms. Specifically, the PG-Schema representation from Example 10 is also a common schema.

**Proposition 1.** For every common schema there exist equivalent SHACL and ShEx schemas.

The translation is relatively straightforward (see Appendix E). The two main observations are that star-free PG-path expressions can be simulated by nested SHACL and ShEx shapes, and that closure of SHACL and ShEx shapes under Boolean connectives allows encoding complex selectors in the shape (as the antecedent of an implication). We illustrate the latter in Example 12.

*Example 12 (Complex paths in selectors).* We want to express that all users who have invited a user who has invited someone (so there is a path following two invited edges) must have a key *email* of type  $\text{str}$ . In PG-schema we express this as:

$$\exists \text{invited} \cdot \text{invited} \Rightarrow \{email : \text{str}\} \& \top$$

At first glance, it seems unclear how to express this in the other formalisms, since they do not permit paths in the selector. However, we can see that paths in selectors can be encoded into the shape: In SHACL, using the same example, we do this by

$$\exists \text{invited} \Rightarrow \neg(\exists \text{invited} \cdot \text{invited}) \vee \exists email.test(\text{str})$$

And in ShEx for this example would be:

$$\{\text{invited}.\{\top\}; \top\} \Rightarrow \neg\varphi_2 \vee \{email.test(\text{str})\}^\circ$$

where  $\varphi_2 = \{(\text{invited}.\varphi_1)^{\geq 1}\}^\circ$  and  $\varphi_1 = \{\text{invited}.\{\top\}^{\geq 1}\}^\circ$ . That is,  $\varphi_1$  is satisfied by nodes that have an outgoing path *invited*, and  $\varphi_2$  by nodes that have an outgoing path *invited* · *invited*. For paths of unbounded length, it is not apparent how such a translation would proceed for ShEx schemas in the absence of recursion.

## 7 RELATED WORK

*SHACL literature.* The authoritative source of SHACL is the W3C recommendation [29]. Further literature on SHACL following its standardisation can be roughly divided into two groups. The first studies the formal properties and expressiveness of the non-recursive fragment [7]. Notable examples in this category (in no particular order) is the work from Delva et al. on data provenance [16], by Pareti et al. on satisfiability and (shape) containment [40] and the work of Leinberger et al. connecting the containment problem to description logics [34]. The other body of work on SHACL is concerned with either proposing a suitable semantics for the recursive fragment [1, 5, 12, 13] or assuming a given one and studying the complexity of certain problems in their chosen recursive setting [39]. First reports on practical applications and use cases for SHACL include, for instance, the expressivity of property constraints, or mining and extracting constraints in the context of practical large KGs such as Wikidata and DBpedia [18, 44].

*ShEx literature.* ShEx was initially proposed in 2014 as a concise and human-readable language to describe, validate, and transform RDF data [43]. Its formal semantics was formally defined in [48]. The semantics of ShEx schemas combining recursion and negation was later presented in [8]. The current semantic specification of the ShEx language has been published as a W3C Community group report [42] and a new language version is currently being defined as part of the IEEE Working group on Shape Expressions<sup>1</sup>. As for practical applications, ShEx has been applied as a descriptive schema

<sup>1</sup><https://shex.io/shex-next/>

language through the Wikidata Schemas project<sup>2</sup>. Additional work went into extending ShEx to handle graph models that go beyond RDF, like WShEx to validate Wikibase graphs [31], ShEx-Star to handle RDF-Star and PShEx to handle property graphs [30].

*PG-Schema literature.* PG-Schema, as introduced in [2], builds on PG-Types and PG-Keys [3] to enhance schema support for property graphs. Despite limited schema support in current systems and the first GQL standard [25], PG-Schema combines flexible type definitions via PG-Types with expressive key constraints from PG-Keys. This formalism provides a robust syntax and semantics for property graph management, aiming to inspire future versions of the GQL standard and broaden the capabilities of graph database systems.

*Works that compare RDF schema formalisms.* In Chapter 7 of [20], the authors compare common features and differences between ShEx and SHACL and [32] presents a simplified language called S, which captures the essence of ShEx and SHACL. Tomaszuk [49] analyzes advancements in RDF validation, highlighting key requirements for validation languages and comparing the strengths and weaknesses of various approaches.

*Interoperability between schema graph formalisms.* Interoperability between schema graph formalisms like RDF and Property Graphs remains challenging due to differences in structure and semantics. RDF focuses on triple-based modeling with formal semantics, while Property Graphs allow flexible annotation of relationships with properties. RDF-star [22] and RDF 1.2 [27] extend RDF 1.1 by enabling statements about triples, aligning more closely with LPG: for instance, RDF-star allows triples to function as subjects or objects, similar to how LPG edges carry properties.

It should be noted here that by adopting *named graphs* [11], already RDF 1.1 provided a mechanism for making statements about (sub-)graphs; likewise, different *reification* mechanisms have been proposed in the literature for RDF in order to “embed” meta-statements about triples (and graphs) in “vanilla” RDF graphs, ranging from the relatively verbose original W3C reification vocabulary as part of the original RDF specification, to more subtle approaches such as singleton property reification [36], which is pretty close to the unique identifiers used for edges in most LPG models. Lastly, custom reification models are used, for instance, in Wikidata, to map Wikibase’s property graph schema to RDF, cf. e.g. [18, 23]. All these approaches, in principle, facilitate general or specific mappings between RDF and LPGs.

Contrary to such mappings using named graphs or reification, there have been several prior proposals for unifying graph data models. The OneGraph initiative [33] also aims to bridge the different graph data models, by promoting a unified graph data model for seamless interaction. Similarly, MilleniumDB’s Domain Graph model [50] aims at providing a general graph model bridging between RDF, RDF-star and property graphs. Work on mappings [4] has also explored schema-independent and schema-dependent methods for transforming RDF into Property Graphs, providing formal foundations for preserving information and semantics.

<sup>2</sup>[https://www.wikidata.org/wiki/Wikidata:WikiProject\\_Schemas](https://www.wikidata.org/wiki/Wikidata:WikiProject_Schemas)

*Schemas for tree-structured data.* The principle of defining (parts of) schemas as a set of pairs  $(sel, \varphi)$  is also used in schema languages for XML. A DTD [10] is essentially such a set of pairs in which  $sel$  selects nodes with a certain label, and  $\varphi$  describes the structure of their children. In XML Schema, the principle was used for defining key constraints (using *selectors* and *fields*) [19, Section 3.11.1]. The equally expressive language BonXai [35] is based on writing the entire schema using such rules. Schematron [24] is another XML schema language that differs from grammar-based languages by defining patterns of assertions using XPath expressions [17]. It excels in specifying constraints across different branches of a document tree, where traditional schema paradigms often fall short. Schematron’s rule-based structure, composed of phases, patterns, rules, and assertions, allows for the validation of documents.

*RDF validation.* Last, but not least, it should be noted that the requirement for (constraining) schema languages—besides ontology languages such as OWL and RDF Schema—in the Semantic Web community is much older than the more recent additions of SHACL and ShEx. Earlier proposals in a similar direction include efforts to add constraint readings of Description Logic axioms to OWL, such as OWL Flight [15] or OWL IC [47]. Another approach is Resource Shapes (ReSh) [45], a vocabulary for specifying RDF shapes. The authors of ReSh recognize that RDF terms originate from various vocabularies, and the ReSh shape defines the integrity constraints that RDF graphs are required to satisfy. Similarly, Description Set Profiles (DSP) [37] and SPARQL Inferencing Notation (SPIN) [28] are notable alternatives. While SHACL, ShEx, and ReSh share declarative, high-level descriptions of RDF graph content, DSP and SPIN offer additional mechanisms for validating and constraining RDF data, each with its own strengths and applications.

## 8 CONCLUSIONS

We provided a formal and comprehensive comparison of the three most prominent schema languages in the Semantic Web and Graph Database communities: SHACL, ShEx, and PG-Schema. Through painstaking discussions within our working group, we managed to (1) agree on a common data model that captures features of both RDF and Property Graphs and (2) extract, for each of the languages, a core that we mutually agree on, which we define formally. Moreover, the definitions of (the cores of) each of the schema languages on a common formal framework allows readers to maximally leverage their understanding of one schema language in order to understand the others. Furthermore, this common framework allowed us to extract the Common Graph Schema Language, which is a cleanly defined set of functionalities shared by SHACL, ShEx, and PG-Schema. This commonality can serve as a basis for future efforts in integrating or translating between the languages, promoting interoperability in applications that rely on heterogeneous data models. For example, we want to investigate recursive ShEx and more expressive query languages for PG-Schema more deeply.

## REFERENCES

- [1] Medina Andresel, Julien Corman, Magdalena Ortiz, Juan L. Reutter, Ognjen Savkovic, and Mantas Simkus. 2020. Stable Model Semantics for Recursive SHACL. In *WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*, Yennun Huang, Irwin King, Tie-Yan Liu, and Maarten van Steen (Eds.). ACM / IW3C2, 1570–1580. <https://doi.org/10.1145/3366423.3380229>



- [2] Renzo Angles, Angela Bonifati, Stefania Dumbrava, George Fletcher, Alastair Green, Jan Hidders, Bei Li, Leonid Libkin, Victor Marsault, Wim Martens, Filip Murlak, Stefan Plantikow, Ognjen Savkovic, Michael Schmidt, Juan Sequeda, Slawek Staworko, Dominik Tomaszuk, Hannes Voigt, Domagoj Vrgoc, Mingxi Wu, and Dusan Zivkovic. 2023. PG-Schema: Schemas for Property Graphs. *Proc. ACM Manag. Data* 1, 2, Article 198 (June 2023), 25 pages. <https://doi.org/10.1145/3589778>
- [3] Renzo Angles, Angela Bonifati, Stefania Dumbrava, George Fletcher, Keith W. Hare, Jan Hidders, Victor E. Lee, Bei Li, Leonid Libkin, Wim Martens, Filip Murlak, Josh Perryman, Ognjen Savkovic, Michael Schmidt, Juan Sequeda, Slawek Staworko, and Dominik Tomaszuk. 2021. PG-Keys: Keys for Property Graphs. In *Proceedings of the 2021 International Conference on Management of Data (Virtual Event, China) (SIGMOD '21)*. Association for Computing Machinery, New York, NY, USA, 2423–2436. <https://doi.org/10.1145/3448016.3457561>
- [4] Renzo Angles, Harsh Thakkar, and Dominik Tomaszuk. 2020. Mapping RDF databases to property graph databases. *IEEE Access* 8 (2020), 86091–86110.
- [5] Bart Bogaerts and Maxime Jakobowski. 2021. Fixpoint Semantics for Recursive SHACL. In *Proceedings 37th International Conference on Logic Programming (Technical Communications), ICLP Technical Communications 2021, Porto (virtual event), 20-27th September 2021 (EPTCS, Vol. 345)*, Andrea Formisano, Yan-hong Annie Liu, Bart Bogaerts, Alex Brik, Verónica Dahl, Carmine Dodaro, Paul Fodor, Gian Luca Pozzato, Joost Vennekens, and Neng-Fa Zhou (Eds.), 41–47. <https://doi.org/10.4204/EPTCS.345.14>
- [6] Bart Bogaerts, Maxime Jakobowski, and Jan Van den Bussche. 2022. SHACL: A Description Logic in Disguise. In *Logic Programming and Nonmonotonic Reasoning - 16th International Conference, LPNMR 2022, Genova, Italy, September 5-9, 2022, Proceedings (Lecture Notes in Computer Science, Vol. 13416)*, Georg Gottlob, Daniela Inlezan, and Marco Maratea (Eds.). Springer, 75–88. [https://doi.org/10.1007/978-3-031-15707-3\\_7](https://doi.org/10.1007/978-3-031-15707-3_7)
- [7] Bart Bogaerts, Maxime Jakobowski, and Jan Van den Bussche. 2024. Expressiveness of SHACL Features and Extensions for Full Equality and Disjointness Tests. *Logical Methods in Computer Science* Volume 20, Issue 1 (Feb. 2024). [https://doi.org/10.46298/lmcs-20\(1:16\)2024](https://doi.org/10.46298/lmcs-20(1:16)2024)
- [8] Iovka Boneva, Jose E. Labra Gayo, and Eric G. Prud'hommeaux. 2017. Semantics and Validation of Shapes Schemas for RDF. In *The Semantic Web - ISWC 2017, Claudia d'Amato, Miriam Fernandez, Valentina Tamma, Freddy Lecue, Philippe Cudré-Mauroux, Juan Sequeda, Christoph Lange, and Jeff Heflin (Eds.)*. Springer International Publishing, Cham, 104–120.
- [9] Angela Bonifati, George H. L. Fletcher, Hannes Voigt, and Nikolay Yakovets. 2018. *Querying Graphs*. Morgan & Claypool Publishers. <https://doi.org/10.2200/S00873ED1V01Y201808DTM051>
- [10] Tim Bray, Jean Paoli, C. M. Sperberg-McQueen, Eve Maler, and François Yergeau. 2008. *Extensible Markup Language (XML) 1.0 (Fifth Edition)*. Technical Report. World Wide Web Consortium.
- [11] Jeremy J. Carroll, Christian Bizer, Pat Hayes, and Patrick Stickler. 2005. Named graphs. *Journal of Web Semantics* 3, 4 (2005), 247–267. <https://doi.org/10.1016/j.websem.2005.09.001> World Wide Web Conference 2005—Semantic Web Track.
- [12] Julien Corman, Fernando Florenzano, Juan L. Reutter, and Ognjen Savkovic. 2019. Validating Shacl Constraints over a SPARQL Endpoint. In *The Semantic Web - ISWC 2019 - 18th International Semantic Web Conference, Auckland, New Zealand, October 26-30, 2019, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 11778)*, Chiara Ghidini, Olaf Hartig, Maria Maleshkova, Vojtech Svátek, Isabel F. Cruz, Aidan Hogan, Jie Song, Maxime Lefrançois, and Fabien Gandon (Eds.). Springer, 145–163. [https://doi.org/10.1007/978-3-030-30793-6\\_9](https://doi.org/10.1007/978-3-030-30793-6_9)
- [13] Julien Corman, Juan L. Reutter, and Ognjen Savkovic. 2018. Semantics and Validation of Recursive SHACL. In *The Semantic Web - ISWC 2018, Denny Vrandečić, Kalina Bontcheva, Mari Carmen Suárez-Figueroa, Valentina Presutti, Irene Celino, Marta Sabou, Lucie-Aimée Kaffee, and Elena Simperl (Eds.)*. Springer International Publishing, Cham, 318–336.
- [14] R. Cyganiak, D. Wood, and M. Lanthaler. 2014. *RDF 1.1 Concepts and Abstract Syntax*. W3C Recommendation. W3C. <http://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/>.
- [15] Jos De Bruijn, Rubén Lara, Axel Polleres, and Dieter Fensel. 2005. OWL DL vs. OWL Flight: Conceptual modeling and reasoning for the Semantic Web. In *Proceedings of the 14th international conference on World Wide Web*. 623–632.
- [16] Thomas Delva, Anastasia Dimou, Maxime Jakobowski, and Jan Van den Bussche. 2023. Data Provenance for SHACL. In *Proceedings 26th International Conference on Extending Database Technology, EDBT 2023, Ioannina, Greece, March 28-31, 2023*. OpenProceedings.org, 285–297. <https://doi.org/10.48786/edbt.2023.23>
- [17] Michael Dyck, Jonathan Robie, and Josh Spiegel. 2017. *XML Path Language (XPath) 3.1*. W3C Recommendation. W3C. <https://www.w3.org/TR/2017/REC-xpath-31-20170321/>.
- [18] Nicolas Ferranti, Jairo Francisco de Souza, Shqiponja Ahmetaj, and Axel Polleres. 2024. Formalizing and Validating Wikidata's Property Constraints using SHACL and SPARQL. *Semantic Web* (2024). <https://doi.org/10.3233/SW-243611>
- [19] Shudi (Sandy) Gao, C. M. Sperberg-McQueen, Henry S. Thompson, Noah Mendelsohn, David Beech, and Murray Maloney. 2012. *W3C XML Schema Definition Language (XSD) 1.1 Part 1: Structures*. Technical Report. World Wide Web Consortium.
- [20] José Emilio Labra Gayo, Eric Prud'hommeaux, Iovka Boneva, and Dimitris Kontokostas. 2017. *Validating RDF Data*. Morgan & Claypool Publishers. <https://doi.org/10.2200/S00786ED1V01Y201707WBE016>
- [21] J. Labra Gayo, H. Knublauch, and D. Kontokostas. 2024. *SHACL Test Suite and Implementation Report*. W3C Document. <https://w3c.github.io/data-shapes/data-shapes-test-suite/>.
- [22] Olaf Hartig. 2014. Reconciliation of RDF\* and Property Graphs. arXiv:1409.3288 [cs.DB] <https://arxiv.org/abs/1409.3288>
- [23] Daniel Hernández, Aidan Hogan, and Markus Krötzsch. 2015. Reifying RDF: What Works Well With Wikidata?. In *Proceedings of the 11th International Workshop on Scalable Semantic Web Knowledge Base Systems co-located with 14th International Semantic Web Conference (ISWC 2015), Bethlehem, PA, USA, October 11, 2015 (CEUR Workshop Proceedings, Vol. 1457)*, Thorsten Liebig and Achille Fokoue (Eds.). CEUR-WS.org, 32–47. [https://ceur-ws.org/Vol-1457/SSWS2015\\_paper3.pdf](https://ceur-ws.org/Vol-1457/SSWS2015_paper3.pdf)
- [24] International Organization for Standardization. 2020. *ISO/IEC 19757-3:2020 Information technology - Document Schema Definition Languages (DSDL) - Part 3: Rule-based validation using Schematron*. Standard. International Organization for Standardization, Geneva, CH.
- [25] International Organization for Standardization. 2024. *ISO/IEC 39075:2024 Information technology - Database languages - GQL*. Standard. International Organization for Standardization, Geneva, CH.
- [26] Maxime Jakobowski. 2024. *Shapes Constraint Language: Formalization, Expressiveness, and Provenance*. Ph.D. Dissertation. Universiteit Hasselt and Vrije Universiteit Brussel.
- [27] Gregg Kellogg, Pierre-Antoine Champin, Olaf Hartig, and Andy Seaborne. 2024. *RDF 1.2 Concepts and Abstract Syntax*. W3C Working Draft. W3C. <https://www.w3.org/TR/2024/WD-rdf12-concepts-20240822/>.
- [28] Holger Knublauch, James A. Hendler, and Kingsley Idehen. 2011. *SPIN - Overview and Motivation*. Technical Report. World Wide Web Consortium.
- [29] H. Knublauch and D. Kontokostas. 2017. *Shapes constraint language (SHACL)*. W3C Recommendation. W3C. <https://www.w3.org/TR/shacl/>.
- [30] José Emilio Labra Gayo. [n. d.]. Extending Shape Expressions for different types of knowledge graphs. In *1st Workshop on Data Quality meets Machine Learning and Knowledge Graphs, DQMLKG, part of Extended Semantic Web Conference 2024, ESWC24 (CEUR Workshop Proceedings, Vol. 3714)*, Sanju Tiwari, Nandana Mihindukulasooriya, Francesco Osborne, Dimitris Kontokostas, Jennifer D'Souza, Mayank Kejriwal, Maria Angela Pellegrino, Anisa Rula, Jose Emilio Labra-Gayo, Michael Cochez, and Mehwish Alam (Eds.). CEUR-WS.org.
- [31] Jose-Emilio Labra-Gayo. 2022. WShEx: A language to describe and validate Wikibase entities. In *Proceedings of the 3rd Wikidata Workshop 2022 co-located with the 21st International Semantic Web Conference (ISWC2022), Virtual Event, Hangzhou, China, October 2022 (CEUR Workshop Proceedings, Vol. 3262)*, Lucie-Aimée Kaffee, Simon Razniewski, Gabriel Amaral, and Kholoud Saad Alghamdi (Eds.). CEUR-WS.org. <https://ceur-ws.org/Vol-3262/paper3.pdf>
- [32] Jose Emilio Labra Gayo, Herminio García-González, Daniel Fernández-Alvarez, and Eric Prud'hommeaux. 2019. Challenges in RDF Validation. In *Current Trends in Semantic Web Technologies: Theory and Practice*, Giner Alor-Hernández, José Luis Sánchez-Cervantes, Alejandro Rodríguez-González, and Rafael Valencia-García (Eds.). Springer, 121–151. [https://doi.org/10.1007/978-3-030-06149-4\\_6](https://doi.org/10.1007/978-3-030-06149-4_6)
- [33] Ora Lassila, Michael Schmidt, Olaf Hartig, Brad Bebe, Dave Bechberger, Willem Broekema, Ankesh Khandelwal, Kelvin Lawrence, Carlos Manuel Lopez-Enriquez, Ronak Sharda, et al. 2023. The OneGraph vision: Challenges of breaking the graph model lock-in 1. *Semantic Web* 14, 1 (2023), 125–134.
- [34] M. Leinberger, P. Seifer, T. Rienstra, R. Lämmel, and S. Staab. 2020. Deciding SHACL Shape Containment through Description Logics Reasoning. In *ISWC'20 (LNCS 12506)*. Springer, 366–383.
- [35] Wim Martens, Frank Neven, Matthias Niewerth, and Thomas Schwentick. 2017. BonXai: Combining the Simplicity of DTD with the Expressiveness of XML Schema. *ACM Trans. Database Syst.* 42, 3 (2017), 15:1–15:42. <https://doi.org/10.1145/3105960>
- [36] Vinh Nguyen, Olivier Bodenreider, and Amit P. Sheth. 2014. Don't like RDF reification?: making statements about statements using singleton property. In *23rd International World Wide Web Conference, WWW '14, Seoul, Republic of Korea, April 7-11, 2014*, Chin-Wan Chung, Andrei Z. Broder, Kyuseok Shim, and Torsten Suel (Eds.). ACM, 759–770. <https://doi.org/10.1145/2566486.2567973>
- [37] Mikael Nilsson. 2008. *Description Set Profiles: A constraint language for Dublin Core Application Profiles*. Technical Report. Dublin Core.
- [38] Paolo Paretì and George Konstantinidis. 2021. A Review of SHACL: From Data Validation to Schema Reasoning for RDF Graphs. In *Reasoning Web. Declarative Artificial Intelligence - 17th International Summer School 2021, Leuven, Belgium, September 8-15, 2021, Tutorial Lectures (Lecture Notes in Computer Science, Vol. 13100)*, Mantas Simkus and Ivan Varzinczak (Eds.). Springer, 115–144. [https://doi.org/10.1007/978-3-030-95481-9\\_6](https://doi.org/10.1007/978-3-030-95481-9_6)
- [39] P. Paretì, G. Konstantinidis, and F. Mogavero. 2022. Satisfiability and Containment of Recursive SHACL. *JWS* 74 (2022), 100721:1–24.

- [40] P. Paret, G. Konstantinidis, F. Mogavero, and T.J. Norman. 2020. SHACL Satisfiability and Containment. In *ISWC'20 (LNCS 12506)*. Springer, 474–493.
- [41] Eric Prud'hommeaux and Thomas Baker. 2017. *ShapeMap Structure and Language*. W3C Draft Community Group Report. W3C. <http://shex.io/shape-map/>.
- [42] Eric Prud'hommeaux, Iovka Boneva, Jose Emilio Labra Gayo, and Gregg Kellog. 2019. *Shape Expressions Language 2.1*. W3C Community Group Report. W3C. <http://shex.io/shex-antics/>.
- [43] Eric Prud'hommeaux, Jose Emilio Labra Gayo, and Harold Solbrig. 2014. Shape expressions: an RDF validation and transformation language. In *Proceedings of the 10th International Conference on Semantic Systems (Leipzig, Germany) (SEM '14)*. Association for Computing Machinery, New York, NY, USA, 32–40. <https://doi.org/10.1145/2660517.2660523>
- [44] Kashif Rabbani, Matteo Lissandrini, and Katja Hose. 2023. Extraction of Validating Shapes from very large Knowledge Graphs. *Proc. VLDB Endow.* 16, 5 (2023), 1023–1032. <https://doi.org/10.14778/3579075.3579078>
- [45] Arthur Ryman. 2014. *Resource Shape 2.0*. Technical Report. World Wide Web Consortium.
- [46] Sherif Sakr, Angela Bonifati, Hannes Voigt, Alexandru Iosup, Khaled Ammar, Renzo Angles, Walid G. Aref, Marcelo Arenas, Maciej Besta, Peter A. Boncz, Khuzaima Daudjee, Emanuele Della Valle, Stefania Dumbrava, Olaf Hartig, Bernhard Haslhofer, Tim Hegeman, Jan Hidders, Katja Hose, Adriana Iammitchi, Vasiliki Kalavri, Hugo Kapp, Wim Martens, M. Tamer Özsu, Eric Peukert, Stefan Plantikow, Mohamed Ragab, Matei Ripeanu, Semih Salihoglu, Christian Schulz, Petra Selmer, Juan F. Sequeda, Joshua Shinavier, Gábor Szárnyas, Riccardo Tommasini, Antonino Tumeo, Alexandru Uta, Ana Lucia Varbanescu, Hsiang-Yun Wu, Nikolay Yakovets, Da Yan, and Eiko Yoneki. 2021. The future is big graphs: a community view on graph processing systems. *Commun. ACM* 64, 9 (2021), 62–71. <https://doi.org/10.1145/3434642>
- [47] Evren Sirin. 2010. Data validation with OWL integrity constraints. In *Web Reasoning and Rule Systems: Fourth International Conference, RR 2010, Bressanone/Brixen, Italy, September 22-24, 2010. Proceedings 4*. Springer, 18–22.
- [48] Slawek Staworko, Iovka Boneva, José Emilio Labra Gayo, Samuel Hym, Eric G. Prud'hommeaux, and Harold R. Solbrig. 2015. Complexity and Expressiveness of ShEx for RDF. In *18th International Conference on Database Theory, ICDT 2015, March 23-27, 2015, Brussels, Belgium*. 195–211. <https://doi.org/10.4230/LIPIcs.ICDT.2015.195>
- [49] Dominik Tomaszuk. 2017. RDF validation: A brief survey. In *Beyond Databases, Architectures and Structures. Towards Efficient Solutions for Data Analysis and Knowledge Representation: 13th International Conference, BDAS 2017, Ustroń, Poland, May 30-June 2, 2017, Proceedings 13*. Springer, 344–355.
- [50] Domagoj Vrgoč, Carlos Rojas, Renzo Angles, Marcelo Arenas, Diego Arroyuelo, Carlos Buil-Aranda, Aidan Hogan, Gonzalo Navarro, Cristian Riveros, and Juan Romero. 2023. MillenniumDB: An Open-Source Graph Database System. *Data Intelligence* (06 2023), 1–39. [https://doi.org/10.1162/dint\\_a\\_00209](https://doi.org/10.1162/dint_a_00209) arXiv:[https://direct.mit.edu/dint/article-pdf/doi/10.1162/dint\\_a\\_00209/2127027/dint\\_a\\_00209.pdf](https://direct.mit.edu/dint/article-pdf/doi/10.1162/dint_a_00209/2127027/dint_a_00209.pdf)

## A DISTILLING THE COMMON DATA MODEL

In this section we discuss the relationship between common graphs and the standard data models of the three schema formalisms—RDF and property graphs.

### A.1 Comparison with RDF

As explained in Section 2, common graphs can be naturally seen as finite sets of triples from  $\mathcal{E} = (\mathcal{N} \times \mathcal{P} \times \mathcal{N}) \cup (\mathcal{N} \times \mathcal{K} \times \mathcal{V})$ , with  $(E, \rho)$  corresponding to  $E \cup \{(u, k, v) \mid \rho(u, k) = v\}$ .

Unlike in RDF, a common graph may contain at most one tuple of the form  $(u, k, v)$  for each  $u \in \mathcal{N}$  and  $k \in \mathcal{K}$ . This reflects the assumption that properties are single-valued, which is present in the property graph data model.

In the RDF context, one would assume the following:

- $\mathcal{N} \subseteq \text{IRIs} \cup \text{Blanks}$ ,
- $\mathcal{P} \subseteq \text{IRIs}$ ,
- $\mathcal{K} \subseteq \text{IRIs}$ ,
- $\mathcal{V} = \text{Literals}$ .

However, the common graph data model does not refer to IRIs, Blanks, and Literals at all, because these are not part of the property graph data model.

In contrast to the RDF model, but in accordance with the perspective commonly taken in databases, both values and nodes are atomic. For nodes we completely abstract away from the actual representation of their identities. We do not even distinguish between IRIs and Blanks. An immediate consequence of this is that schemas do not have access to any information about the node other than the triples in which it participates. In particular, they cannot compare nodes with constants. This is a significant restriction with respect to the RDF data model, but it follows immediately from the same assumption made in the property graph data model. On the positive side, this aspect is entirely orthogonal to the main discussion in this paper, so eliminating it from the common data model does not oversimplify the picture.

For values we take a more subtle approach: we assume a set  $\mathcal{T}$  of value types, with each  $v \in \mathcal{T}$  representing a set  $\llbracket v \rrbracket \subseteq \mathcal{V}$ . This captures uniformly data types, such as integer or string, and user-defined checks, such as interval bounds for numeric values or regular expressions for strings. On the other hand, the common graph data model does not include any binary relations over values, such as an order.

### A.2 Comparison with property graphs

Let us recall the standard definition of property graphs [2].

DEFINITION 14 (PROPERTY GRAPH). A property graph is a tuple  $(N, E, \pi, \lambda, \rho)$  such that

- $N$  is a finite set of nodes;
- $E$  is a finite set of edges, disjoint from  $N$ ;
- $\pi : E \rightarrow (N \times N)$  maps edges to their source and target;
- $\lambda : (N \cup E) \rightarrow 2^{\mathcal{P}}$  maps nodes and edges to finite sets of labels;
- $\rho : (N \cup E) \times \mathcal{K} \rightarrow \mathcal{V}$  is a finite-domain partial function mapping element-key pairs to values.

A common graph  $G = (E', \rho')$  can be easily represented as a property graph by letting

- $N = \text{Nodes}(G)$ ,
- $E = E'$ ,
- $\pi = \{(e, (v_1, v_2)) \mid e = (v_1, p, v_2) \in E\}$ ,
- $\lambda = \{(e, \{p\}) \mid e = (v_1, p, v_2) \in E\} \cup \{(v, \emptyset) \mid v \in N\}$ , and
- $\rho = \rho'$ .

It is possible to characterise exactly the property graphs that are such representations of common graphs. These are the property graphs  $(N, E, \pi, \lambda, \rho)$  for which it holds that:

- (1)  $\lambda(v) = \emptyset$  for all  $v \in N$ , and  $\lambda(e)$  is a singleton for all  $e \in E$ ,
- (2) there cannot be two distinct edges  $e_1, e_2 \in E$  such that  $\pi(e_1) = \pi(e_2)$  and  $\lambda(e_1) = \lambda(e_2)$ , and
- (3)  $\rho(e, k)$  is undefined for all  $e \in E, k \in \mathcal{K}$ .

So, common graphs can be interpreted as restricted property graphs: no labels on nodes, single labels on edges, no parallel edges with the same label, and no properties on edges. All these restrictions are direct consequences of the nature of the RDF data model.

While these restrictions seem severe at a first glance, the resulting data model can actually easily simulate unrestricted property graphs: labels on nodes can be simulated with the presence of corresponding keys, edges can be materialised as nodes if we need properties over edges or parallel edges with the same label. This means not only that common graphs can be used without loss of generality in expressiveness and complexity studies, but also that the corresponding restricted property graphs are flexible enough to be usable in practice, while additionally guaranteeing interoperability with the RDF data model.

### A.3 Class information

The common graph data model does not have direct support for class information. The reason for this is that RDF and property graphs handle class information rather differently. In RDF, both class and instance information is part of the graph data itself: classes are elements of the graph, subclass-superclass relationships are represented as edges between classes, and membership relationships are represented as edges between elements and classes. In property graphs, membership of a node in a class is indicated by a label put on the node. A node can belong to many classes, but the only way to say that class  $A$  is a subclass of class  $B$  is to ensure in the schema that each node with label  $A$  also has label  $B$ . That is,

- in property graphs class membership information is available locally in a node, but consistency must be ensured by the schema,
- in RDF, obtaining class membership information requires navigating in the graph, but consistency is for free.

Clearly, both approaches have their merits, but when passing from one to the other data needs to be translated. This means that we cannot pick one of these approaches for the common data model while keeping it a natural submodel of both RDF and property graphs. Therefore, to reduce the complexity of this study, we do not include in our common data model any dedicated features for supporting class information. Note, however, that common graphs can support both these approaches indirectly: designated predicates can be used to represent membership and subclass relationships, and keys with a dummy value can simulate node labels.

## B STANDARD SHACL

Standard SHACL defines shapes as a conjunction of *constraint components*. The different constructs from our formalization correspond to fundamental building blocks of these constraint components. Next to that, the formalization of SHACL presented in this paper is less expressive than standard SHACL: both because we defined it here for the common data model (which is a strict subset of RDF), and because we want to simplify our narrative. For example, we leave out the comparison of RDF terms using `sh:lessThan` to simplify our story. Furthermore, because the common graph data-model does not define language tags, the corresponding constraint components from standard SHACL are left out.

Our formalization is closely tied to the ones found in the literature. There, correspondence between the formalization of the literature and standard SHACL has been shown in detail [26]. This section highlights and discusses some relevant details.

*Class targets and constraint component.* As a consequence of the common graph data model not directly supporting the modelling of classes, some class-based features are not adapted in our formalization. Specifically, there are no selectors (“target declarations” in standard SHACL) that involve classes. Furthermore, the value type constraint component `sh:class` is not covered by our formalization.

*Closedness.* In standard SHACL syntax, closedness is a property that takes a true or false value. The semantics of closedness is based on a list of predicates that are allowed for a given focus node. This list can be inferred based on the predicated used in property shapes, or this list can be explicitly given using the `sh:ignoredProperties` keyword. Our formalization effectively adopts the latter approach: `closed(Q)` means that the properties mentioned in the set  $Q$  are the ignored properties.

*Path expressions.* The path expressions used in our formalization deviate from the standard in three obvious ways. First, we make a distinction between ‘keys’ and ‘predicates’. This is simply a consequence of using our common graph data model. Second, we leave out some of the immediately available path constructs from standard SHACL: one-or-more paths and zero-or-one path. However, these are expressible using the building blocks of our formalization: one-or-more paths are expressed as  $\pi \cdot \pi^*$ , and zero-or-one paths are expressed as  $\pi \cup \text{id}$ . Lastly, our path expression allow for writing the identity relation explicitly. This cannot be done in literal standard SHACL syntax, but its addition to the formalization serves to highlight its hidden presence in the language. Writing the identity relation directly in a counting construct, e.g.,  $\exists^{\geq n} \text{id} \cdot \top$ , never adds expressive power. In the case of  $n = 1$ , the shape is always satisfied (and thus equivalent to  $\top$ ), and it is easy to see that for any  $n > 1$ , it is never satisfied. The situation with complex path expressions in counting constructs is less clear from the outset. However, it has been shown [7] (Lemma 3.3), that the only case where `id` adds expressiveness is with complex path expressions of the form  $\pi \cup \text{id}$ . This is exactly the definition of zero-or-one paths and is therefore covered by standard SHACL. Another place where `id` can occur in our formalization is in the equality and disjointness constraints, e.g., `eq(id, p)`. According to the standard SHACL recommendation,

you cannot write this shape. However, in the SHACL Test Suite [21] test `core/node/equals-001`, the following shape is tested for:

```
ex:TestShape
  rdf:type sh:NodeShape ;
  sh:equals ex:property ;
  sh:targetNode ex:ValidResource1 .
```

on the following data:

```
ex:ValidResource1
  ex:property ex:ValidResource1 .
```

The intended meaning of this test is, in natural language: “The `targetNode ex:ValidResource1` has an `ex:property` self-loop and no other `ex:property` properties”. Effectively, this is the semantics for our `eq(id, p)` construct. The situation with `disj(id, p)` is similar.

We therefore have an ambiguous situation: the standard description of SHACL does not allow for shapes of the form `eq(id, p)`, but the test suite, and therefore all implementations that pass it completely, do<sup>3</sup>. It then seems fair to include this powerful construct in the formalization.

*Comparisons with constants.* A direct consequence of the assumption that node identities in the common graph data model are hidden from the user, our abstraction of SHACL on common graphs does not support comparisons with constants from  $\mathcal{N}$ . Comparisons with constants from  $\mathcal{V}$  are allowed.

*Node tests.* Our formalization uses the `test( $\forall$ )` construct to denote many of the node tests available in SHACL. We list the tests from standard SHACL that are covered by this construct.

**DatatypeConstraintComponent** Tests whether a node has a certain datatype.

**MinExclusiveConstraintComponent**

**MinInclusiveConstraintComponent**

**MaxExclusiveConstraintComponent**

**MaxInclusiveConstraintComponent** These four constraints cover can check whether a node is larger (**Max**) or smaller (**Min**) than some value, and whether this forms a partial order (**Inclusive**) or a strict, or total, order (**Exclusive**). Based on the SPARQL `<` or `≤` operator mapping.

**MaxLengthConstraintComponent**

**MinLengthConstraintComponent** These two constraints test whether the length of the lexical form of the node is “larger” or equal (resp. “smaller” or equal) than some provided integer value. Strictly speaking, the recommendation defines these constraint components also on IRIs. However, we limit their use to Literals.

**PatternConstraintComponent** Tests whether the length of the lexical form of the node satisfies some regular expression. Strictly speaking, the recommendation defines these constraint components also on IRIs. However, we limit their use to Literals.

Then there are two types of tests not covered by our formalization:

<sup>3</sup>Incidentally, all implementations currently mentioned in the implementation report handle these cases correctly.

**NodeKindConstraintComponent** Tests whether a node is an IRI, Blank Node, or Literal. Our tests apply only to RDF Literals.

**LanguageInConstraintComponent** Test whether the language tag of the node is one of the specified language tags. This feature is not supported by our data model, since it lacks language tags.

## C STANDARD SHEX

The Shape Expressions Language (ShEx) [42] and the ShapeMaps language [41] are defined by the Shape Expressions Community group<sup>4</sup> at W3C. Hereafter we use *standard ShEx* to refer to the language defined in [42], while *ShEx* designates the language presented in this paper.

In this appendix we support the following

**Claim 1.** *On common graphs, the expressive power of ShEx schemas is equivalent to the expressive power of standard ShEx schemas without recursion.*

Section C.3 explains standard ShEx on common graphs, while Section C.2 explains standard ShEx without recursion.

### C.1 Standard ShEx schema and the validation problem

A standard ShEx schema is a set of named shape expressions, and is usually formalised as a pair  $(L, def)$  where  $L$  is a finite set of shape names (in practice these are IRIs), and  $def$  is a function that associates a shape expression with every shape name. In standard ShEx, the validation problem  $\mathcal{G} \models \mathcal{S}$  from Section 2.4 is not defined as such. That is, the ShEx specification [42] does not say what it means for a graph to be valid w.r.t. a standard ShEx schema; it only defines what it means for a node in a graph to satisfy a shape expression.

However, the problem considered in practice is whether some selected nodes in the graph satisfy some indicated shape expressions from the schema. This is specified by a shape map [41]. A shape map can be formalised as a set of pairs of the form  $(sel, l)$ , where  $l \in L$  and  $sel$  is a unary query. While the shape map specification [41] allows the selectors from Definition 7, most implementations allow general SPARQL queries as selectors.

In the current paper, we integrate the shape map into the schema itself, which allows us to specify the validation problem in a uniform way for the three graph schema formalisms considered. Additionally, in shape maps we do not use shape names but shape expressions directly; the next section argues why this is not a problem from the point of view of expressive power.

### C.2 Shape names and recursion

Recursion is an important mechanism in standard ShEx. In the current paper however we consider only ShEx schemas without recursion (to be defined shortly). This restriction was made because neither standard SHACL<sup>5</sup> nor PG-Schema allow for recursion.

<sup>4</sup><https://www.w3.org/community/shex/>

<sup>5</sup>There exist formalizations of SHACL that introduce recursion, but recursion is not specified in the standard.

In standard ShEx, the fact that shape expressions are named allows to refer to them using their name. In particular, references allow for circular recursive definitions. As an example, consider the standard ShEx schema on Figure 2. It contains the single shape name `ex:User` whose definition is given by the shape expression inside the curly braces. The latter shape expression refers to itself: `@ex:User` indicates a reference to the shape expression named `ex:User`. Concretely, the shape expression requires from an RDF node to have an `ex:email` predicate whose value is a string, as well as any number of `ex:invited` predicates whose values are nodes that satisfy the shape expression named `ex:User`.

```

PREFIX ex: <http://ex.example/#>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
ex:User {
  ex:email xsd:string ;
  ex:invited @ex:User *
}

```

Figure 2: A standard ShEx schema.

A standard ShEx schema  $(L, def)$  is considered with recursion when there is a shape name  $l \in L$  whose definition uses a reference `@` to itself, either directly or transitively through references to other shape names. Every standard ShEx schema without recursion can be rewritten to an equivalent schema without references, simply by replacing every reference by its definition. In other words, references in standard ShEx do not add expressive power for non-recursive schemas.

### C.3 Syntax of standard ShEx on common graphs

The syntax and semantics of standard ShEx differ slightly from the ones presented here. There are some purely syntactic differences, as for instance the boolean operators that are called `and`, `or` and `not` in standard ShEx, for which we adopt here the usual mathematical notation  $\wedge$ ,  $\vee$  and  $\neg$ . Another difference concerns the so-called node constraints of standard ShEx. These are constraints to be verified on the actual node of an RDF graph (which is an IRI, a literal or a blank node) without considering its neighbourhood. As pointed out in Section A.1, node constraints for nodes (i.e. elements of  $\mathcal{N}$ ) are irrelevant for common graphs. Standard ShEx node constraints on values correspond to the atomic shape expressions `test( $\forall$ )` and `is( $c$ )` in ShEx. Their expressive power can be entirely captured by selecting for  $\mathcal{T}$  a language equivalent to node constraints on values in standard ShEx. Note finally that the `test( $\text{any}$ )` in ShEx allows to distinguish nodes from values.

The grammar presented on Figure 3 gives an abstract syntax for standard ShEx *shape expressions*  $se$  and *triple expressions*  $te$  for common graphs. Note that we use the ShEx syntax for the boolean operators and for node constraints. The other differences are:

- In standard ShEx, the atomic shape expression that defines the neighbourhood of a node (non-terminal  $sh$ ) is parameterised by a set of extra (possibly inversed) predicates and keys  $Q$ , indicated by the keyword `extra`. In Section C.4.1 we show that extra is syntactic sugar in standard ShEx.
- In standard ShEx, the atomic shape expression derived from the non-terminal  $sh$  can have an optional closed modifier.

On the other hand, ShEx introduces the triple expressions  $\neg P$  and  $\neg P^-$ . As we will see, the latter are used when translating standard ShEx to ShEx in order to distinguish between closed and non-closed standard ShEx shape expressions.

- Triple expressions in ShEx contain the atomic expression  $\varepsilon$ , while standard ShEx does not allow it directly. On the other hand standard ShEx allows to use intervals of the form  $[min; max]$  to define bounded or unbounded repetition, while ShEx allows only the unbounded repetition  $*$ . We show in Section C.4.2 that the two variants have equivalent expressive power.
- Triple constraints in standard ShEx (rule  $tc$ ) allow to use a `.` (dot) instead of the shape expression, which is in fact equivalent to the ShEx shape expression  $\{\top\}$ .

Note that, strictly speaking, the extra predicates and keys are optional in standard ShEx. However, an absent extra set is equivalent to extra  $\emptyset$ , therefore we will consider that it is always present.

$$\begin{aligned}
 se & ::= is(c) \mid test(\forall) \mid sh \mid closed\ sh \mid se \wedge se \mid se \vee se \mid \neg se . \\
 sh & ::= extra\ Q\ \{te\} . \\
 te & ::= tc \mid te ; te \mid te|te \mid te[min; max] . \\
 tc & ::= q\ se \mid q . .
 \end{aligned}$$

with  $c \in \mathcal{N} \cup \mathcal{V}$ ,  $\forall \in \mathcal{T}$ ,  $q \in \mathcal{P} \cup \mathcal{K} \cup \mathcal{P}^- \cup \mathcal{K}^-$ ,  
 $Q \subseteq_{fin} \mathcal{P} \cup \mathcal{K} \cup \mathcal{P}^- \cup \mathcal{K}^-$ ,  $min \in \mathbb{N}$  and  $max \in \mathbb{N} \cup \{*\}$ .

Figure 3: Abstract syntax for standard ShEx.

### C.4 Translations between ShEx and standard ShEx

In this section we introduce a back and forth translation between standard ShEx without recursion and ShEx. We claim that these translations preserve the semantics w.r.t. the validity of a graph. The claim is presented without proof as it would require to define formal semantics for standard ShEx. The proof is however not difficult to make using the formal semantics from [8].

*C.4.1 Eliminating extra from standard ShEx.* We show on an example how the extra constructed can be eliminated in standard ShEx. The same example will be used later on for the translation from standard ShEx to ShEx, therefore it is described in detail.

*Example 13.* Consider the standard ShEx shape expression

$$\begin{aligned}
 se & = extra\ \{p_1, p_2\}\ \{te\} \\
 \text{with} \quad te & = p_1\ \{p.\} ; p_1\ \{p'.\} ; p_3.\ ; p_4^- . \\
 \text{and} \quad p_1, p_2, p_3, p_4, p, p' & \in \mathcal{P} \cup \mathcal{K}
 \end{aligned}$$

that has a set of extra predicates and keys  $\{p_1, p_2\}$ . It is satisfied by nodes whose neighbourhood has the following outgoing triples:

- (1) one  $p_1$ -triple leading to a node that satisfies  $\{p.\}$ ,
- (2) another  $p_1$ -triple leading to a node that satisfies  $\{p'.\}$ ,
- (3) a  $p_3$ -triple leading to an unconstrained node,

- (4) because  $p_1$  appears as extra, other  $p_1$ -triples are also allowed as long as they satisfy **none** of the constraints present for  $p_1$  in  $te$ , that is, they satisfy neither  $\{p \cdot\}$  nor  $\{p' \cdot\}$ ,
- (5) because  $p_2$  appears as extra,  $p_2$ -triples are allowed and their target is not constrained because  $p_2$  does not appear in the triple expression  $te$ ,
- (6) finally, because the shape is not closed, all outgoing triples whose predicate is not  $p_1$  neither  $p_3$  are allowed,  $\{p_1, p_3\}$  being the set of non-inversed predicates that occur directly in  $te$ . Formally,  $\{p_1, p_3\} = \text{preds}(te) \cap \mathcal{P} \cap \mathcal{K}$ , see Section C.4.3.

The node has also the following incoming triples:

- (7) one incoming  $p_4$ -triple coming from an unconstrained node,
- (8) because standard ShEx does not close the incoming triples, all incoming triples whose predicate is different from  $p_4$  are allowed.

The shape expression  $se$  from Example 13 is equivalent to the following shape expression without extra:

$$\{te ; te_{p_1}^* ; te_{p_2}^*\}$$

where

$$te_{p_1} = p_1 (\neg \{p \cdot\} \wedge \neg \{p' \cdot\}) \quad \text{and} \quad te_{p_2} = (p_2 \cdot)$$

The sub-expression  $te_{p_1}^*$  allows to satisfy the requirement (4) from Example 13, while the sub-expression  $te_{p_2}^*$  allows to satisfy the requirement (5).

This construction for eliminating extra can be generalised to arbitrary shape expressions. The idea is to combine (with the ; operator) the initial triple expression with a sub-expression of the form  $q se_q^*$  for every (possibly inverse) extra predicate  $q$ , where  $se_q$  is the conjunction of the negated shape expressions  $se'$  such that  $q se'$  directly in  $te$ .

Without loss of generality, **from now on, we consider only standard ShEx shape expressions without extra**.

**C.4.2 Normalised triple expressions.** We now show how standard ShEx triple expressions and triple expressions can be normalised so that they use a limited number of operators, which will be useful for the translation between standard ShEx and ShEx.

*Normalisation of standard ShEx triple expressions.* A standard ShEx triple expression is normalised when it uses only the intervals  $[0; 1]$  and  $[0; *]$ . Standard ShEx triple expressions are normalised using rewriting based on the following equivalences:

$$te[ \min; * ] = te[0; *] ; \underbrace{te ; \dots ; te}_{\min \text{ times}}$$

$$te[ \min; \max ] = \underbrace{te ; \dots ; te}_{\min \text{ times}} ; \underbrace{te[0; 1] ; \dots ; te[0; 1]}_{\max - \min \text{ times}} \quad \text{when } \max \neq *$$

*Normalisation of ShEx triple expressions.* Here after,  $f$  designates a ShEx triple expression derivable from the third rule of the grammar in Definition 6. For every ShEx triple expression  $f$ , we define  $f^? = f | \varepsilon$ . A ShEx triple expression  $f$  is normalised if either  $f = \varepsilon$ , or  $f$  does not use  $\varepsilon$  as sub-expression, but can use the ? operator defined here-above. Every triple expression can be normalised by

eliminating occurrences of  $\varepsilon$  using these two properties and the ? operator:

- it is a neutral element for the ; operator, that is,  $f ; \varepsilon = \varepsilon ; f = f$  for every ShEx triple expression  $f$ ,
- $\varepsilon^* = \varepsilon$ .

Without loss of generality, **from now on we consider only normalised triple expressions**.

**C.4.3 Direct predicates of triple expressions.** This section is devoted to two technical definitions. For every triple expression we define the set of (possibly inversed) predicates and keys that appear directly in the expression. Formally, if  $f$  is a ShEx triple expression derived by the third rule of the grammar in Definition 6, then we define the set  $\text{preds}(f) \subseteq \mathcal{P} \cup \mathcal{K} \cup \mathcal{P}^- \cup \mathcal{K}^-$  inductively on the structure of  $f$  by:

$$\begin{aligned} \text{preds}(\varepsilon) &= \emptyset \\ \text{preds}(p \cdot \varphi) &= \{p\} \\ \text{preds}(p^- \cdot \varphi) &= \{p^-\} \\ \text{preds}(\varphi ; \varphi') &= \text{preds}(\varphi) \cup \text{preds}(\varphi') \\ \text{preds}(\varphi | \varphi') &= \text{preds}(\varphi) \cup \text{preds}(\varphi') \\ \text{preds}(\varphi^*) &= \text{preds}(\varphi) \end{aligned}$$

For a standard ShEx triple expression  $te$ , the set  $\text{preds}(te)$  is defined similarly (recall that  $q \in \mathcal{P} \cup \mathcal{K} \cup \mathcal{P}^- \cup \mathcal{K}^-$ ):

$$\begin{aligned} \text{preds}(q se) &= \{q\} \\ \text{preds}(q \cdot) &= \{q\} \\ \text{preds}(se ; se') &= \text{preds}(se) \cup \text{preds}(se') \\ \text{preds}(se | se') &= \text{preds}(se) \cup \text{preds}(se') \\ \text{preds}(te[\min; \max]) &= \text{preds}(se) \end{aligned}$$

**C.4.4 Translation from standard ShEx to ShEx.** With every standard ShEx shape expression  $se$  we associate the ShEx shape expression  $\tau(se)$  as defined in Table 9. It is defined by mutual recursion with the corresponding translation function  $\tau_e(te)$  for standard ShEx triple expressions  $te$  presented in Table 8.

**Table 8: Translation from standard ShEx to ShEx for normalised triple expressions, with  $q \in \mathcal{P} \cup \mathcal{K} \cup \mathcal{P}^- \cup \mathcal{K}^-$ .**

$te$	$\tau_e(te)$
$q se$	$p. \tau(se)$
$q \cdot$	$p. \{\top\}$
$te ; te'$	$\tau_e(te) ; \tau_e(te')$
$te   te'$	$\tau_e(te)   \tau_e(te')$
$te[0; *]$	$\tau_e(te)^*$
$te[0; 1]$	$\tau_e(te)   \varepsilon$

**C.4.5 Translation from ShEx to standard ShEx.** Unless otherwise specified, in the sequel  $f$  designates a ShEx triple expression produced by the third rule in the grammar in Definition 6. In Table 10 we present a function that with every normalised ShEx triple expression  $f$  associates the standard ShEx triple expression  $\sigma_e(f)$ . It is defined by mutual recursion with the translation function that

**Table 9: Translation from standard ShEx to ShEx for shape expressions.**

$se$	$\tau(se)$
$is(c)$	$is(c)$
$\tau(\text{test}(\mathbb{V}))$	$\text{test}(\mathbb{V})$
$se \wedge se'$	$\tau(se) \wedge \tau(se')$
$se \vee se'$	$\tau(se) \vee \tau(se')$
$\neg se$	$\neg \tau(se)$
closed $\{te\}$	$\tau_e(te); (\neg Q^-)^*$ with $Q = \text{preds}(te) \cap (\mathcal{P}^- \cup \mathcal{K}^-)$
$\{te\}$	$\tau_e(te); (\neg Q^-)^*; (\neg P)^*$ with $Q = \text{preds}(te) \cap (\mathcal{P}^- \cup \mathcal{K}^-)$ and $P = \text{preds}(te) \cap (\mathcal{P} \cup \mathcal{K})$

with every ShEx shape expression  $\varphi$  associates a standard ShEx shape expression  $\sigma(\varphi)$ , and that will be presented shortly. Note that the case  $f = \varepsilon$  is omitted in Table 10: recall that in normalised ShEx triple expressions,  $\varepsilon$  can only appear standalone (not in sub-expressions), therefore this case  $f = \varepsilon$  will be treated with shape expressions.

**Table 10: Translation from ShEx to standard ShEx for normalised triple expressions.**

$f$	$\sigma_e(f)$
$p.\varphi$	$p \sigma_e(\varphi)$
$p^-.\varphi$	$p^- \sigma_e(\varphi)$
$f; f'$	$\sigma_e(f); \sigma_e(f')$
$f   f'$	$\sigma_e(f)   \sigma_e(f')$
$f^*$	$\sigma_e(f)[0; *]$
$f^?$	$\sigma_e(f)[0; 1]$

The definition of  $\sigma(\varphi)$  is straightforward for the following cases:

$$\begin{aligned} \sigma(is(c)) &= is(c) \\ \sigma(\text{test}(\mathbb{V})) &= \text{test}(\mathbb{V}) \\ \sigma(\varphi \wedge \varphi') &= \sigma(\varphi) \wedge \sigma(\varphi') \\ \sigma(\varphi \vee \varphi') &= \sigma(\varphi) \vee \sigma(\varphi') \\ \sigma(\neg \varphi) &= \neg \sigma(\varphi) \end{aligned}$$

The remaining case is for a shape expression of the form  $\{e\} = \{f; \dots\}$  where  $f$  is normalised. Consider the most general case

$$e = f; (\neg P^-)^*; (\neg Q)^*$$

Let also

$$\begin{aligned} \{p_1, \dots, p_m\} &= (\text{preds}(f) \cap \mathcal{P}^- \cap \mathcal{K}^-) \setminus P \\ \{q_1, \dots, q_n\} &= (\text{preds}(f) \cap \mathcal{P} \cap \mathcal{K}) \setminus Q. \end{aligned}$$

Intuitively,  $\{p_1, \dots, p_m\}$  is the set of predicates that are not allowed to appear on incoming edges in the neighbourhoods defined by  $e$ , and similarly  $\{q_1, \dots, q_n\}$  are the forbidden predicates for outgoing

edges. Then

$$\sigma(\{e\}) = \left\{ \begin{array}{l} \sigma_e(f); \\ p_1^-.[0; 0]; \dots; p_m^-.[0; 0]; \\ q_1.[0; 0]; \dots; q_n.[0; 0] \end{array} \right\}$$

If  $f = \varepsilon$ , then the term  $\sigma_e(f)$  on the first line of the definition of  $\sigma(\{e\})$  must be omitted.

The remaining case for the definition of  $\sigma(\{e\})$  is for

$$e = f; (\neg P^-)^*$$

Let  $\{p_1, \dots, p_m\}$  be as before. Then

$$\sigma(\{e\}) = \text{closed} \left\{ \begin{array}{l} \sigma_e(f); \\ p_1^-.[0; 0]; \dots; p_m^-.[0; 0] \end{array} \right\}$$

As before, if  $f = \varepsilon$ , then the term  $\sigma_e(f)$  must be omitted.

## D STANDARD PG-SCHEMA

We present here a version of PG-Schema that is restricted to common graphs, and therefore as such a simplified version of the original PG-Schema, but stays closer to the original form of PG-Schema as discussed in [2] than the version in the body of this paper. It is meant to illustrate how the version of PG-Schema in the body describes indeed a reasonable core of the original full PG-Schema.

The central idea of standard PG-Schema and of PG-Schema in [2] is that a schema (called *graph type* in this context) consists of three parts: (1) a set of node types, (2) a set of edge types, and (3) a set of graph constraints that represents logical statements about the property graph that must hold for it to be valid. A particular property graph is then said to be valid wrt. such a graph type if (1) every node in the property graph is in the semantics of at least one node type, (2) every edge in the property graph is in the semantics of at least one edge type, and (3) the property graph satisfies all specified graph constraints.

The organisation of this section is as follows. We first discuss the notions of *node types* and *edge types*. After that we discuss how path expressions are defined, after which we discuss what graph constraints look like in this setting. In the final two subsections we discuss how this version of PG-Schema relates the original defined in [2], and how it relates to the one defined in this paper.

### D.1 Node types

The purpose of node types in PG-Schema is to describe nodes, their properties and their labels. Since in the common graph model nodes no longer have labels, node types become simply record types where the record fields describe the allowed keys. In addition it is indicated with these record type whether they are *closed* or *open*, where the first indicates that only the indicated keys are allowed, and the later that additional keys are allowed. This leads to the following formal definition.

**DEFINITION 15 (NODE TYPE).** A node type is an expression  $\mathfrak{n}$  of the form defined by the grammar

$$\mathfrak{n} ::= \{\} \mid \{\}^\circ \mid \{k : \mathbb{V}\} \mid \{k : \mathbb{V}\}^\circ \mid \mathfrak{n} \& \mathfrak{n} \mid \mathfrak{n} \mid \mathfrak{n}.$$

where  $k \in \mathcal{K}$  and  $\mathbb{V} \in \mathcal{T}$ .

Here  $\{\}^\circ$  indicates the open record type without any keys,  $\{k : \mathbb{V}\}^\circ$  is the open record type with just key  $k$  of type  $\mathbb{V}$ .

**Table 11: Semantics of standard PG-path expressions.**

$\pi$	$\llbracket \pi \rrbracket^{\mathcal{G}} \subseteq (\mathcal{N} \cup \mathcal{V}) \times (\mathcal{N} \cup \mathcal{V})$ for $\mathcal{G} = (E, \rho)$
$\mathfrak{e}$	$\{(u, v) \mid \exists p : (u, p, v) \in E \wedge (\rho(u), p, \rho(v)) \in \llbracket \mathfrak{e} \rrbracket\}$
$\neg \mathfrak{e}$	$\{(u, v) \mid \exists p : (u, p, v) \in E \wedge (\rho(u), p, \rho(v)) \notin \llbracket \mathfrak{e} \rrbracket\}$

For these node types  $\mathfrak{m}$  we define a value semantics  $\llbracket \mathfrak{m} \rrbracket$ , just like for content types, which is defined in the same way except that  $\llbracket \{\}^\circ \rrbracket = \llbracket \top \rrbracket$  and  $\llbracket \{k : \mathfrak{v}\}^\circ \rrbracket = \llbracket \{k : \mathfrak{v}\} \& \top \rrbracket$ . It follows straightforwardly that for each content type there is an equivalent node type and vice versa.

## D.2 Edge types

In PG-Schema there is a notion of edge type, which consists of three parts: (1) a type describing the source node, (2) a type describing describing the contents of the edge itself, and (3) a type describing the the target node. Since in common graphs the content of an edge is just a label, a type describing this content can be simply an expression of the form  $\star$  (a wild-card indicating that any label is possible) or a finite set  $P$  of labels (indicating that only these labels are allowed). So we get the following definition for edge types.

**DEFINITION 16 (EDGE TYPE).** *An edge type is an expression  $\mathfrak{e}$  of the form defined by the grammar*

$$\mathfrak{e} ::= \mathfrak{m} \xrightarrow{\star} \mathfrak{m} \mid \mathfrak{m} \xrightarrow{P} \mathfrak{m} \mid \mathfrak{e} \& \mathfrak{e} \mid \mathfrak{e} \mid \mathfrak{e} .$$

where  $P$  is a finite subset of  $\mathcal{P}$ .

As for node types, we define for edge types a value semantics, which in this case defines which combinations of (1) source node content, (2) edge content, and (3) target node content are allowed.

**DEFINITION 17 (VALUE SEMANTICS OF EDGE TYPES).** *With an edge type  $\mathfrak{e}$  we associate a value semantics  $\llbracket \mathfrak{e} \rrbracket \subseteq \mathcal{R} \times \mathcal{P} \times \mathcal{R}$  which is defined with induction on the structure of  $\mathfrak{e}$  as follows:*

- (1)  $\llbracket \mathfrak{m}_1 \xrightarrow{\star} \mathfrak{m}_2 \rrbracket = \llbracket \mathfrak{m}_1 \rrbracket \times \mathcal{P} \times \llbracket \mathfrak{m}_2 \rrbracket$
- (2)  $\llbracket \mathfrak{m}_1 \xrightarrow{P} \mathfrak{m}_2 \rrbracket = \llbracket \mathfrak{m}_1 \rrbracket \times P \times \llbracket \mathfrak{m}_2 \rrbracket$
- (3)  $\llbracket \mathfrak{e}_1 \& \mathfrak{e}_2 \rrbracket = \{ ((r_1 \cup s_1), p, (r_2 \cup s_2)) \in \mathcal{R} \times \mathcal{P} \times \mathcal{R} \mid (r_1, p, r_2) \in \llbracket \mathfrak{e}_1 \rrbracket \wedge (s_1, p, s_2) \in \llbracket \mathfrak{e}_2 \rrbracket \}$
- (4)  $\llbracket \mathfrak{e}_1 \mid \mathfrak{e}_2 \rrbracket = \llbracket \mathfrak{e}_1 \rrbracket \cup \llbracket \mathfrak{e}_2 \rrbracket$

## D.3 Path expressions

We define here a notion of path expression that we call *standard PG-path expression* and that is similar to the notion of PG-path expression of Definition 9, except that in the positions where a content type  $\mathfrak{c}$  is allowed, we allow the use of a node type  $\mathfrak{m}$  or edge type  $\mathfrak{e}$ . As discussed earlier, node types are equivalent to content types, and so we can given them the same semantics in path expressions. Edge types, however, are different, and their semantics in path expressions is given in Table 11.

## D.4 Graph constraints

The graph constraints in PG-Schema are based on the constraints discussed in PG-Keys [3]. Although the latter paper focuses on key constraints, it also discusses other closely related cardinality

constraints. We capture these constraints here in the context of the common graph data model with the following formal definition.

**DEFINITION 18 (PG-CONSTRAINT).** *A PG-constraint is a formula of one of the following three forms:*

- K1:**  $\forall x : \varphi(x) \Leftarrow \mathbf{Key} \bar{y} : \psi(x, \bar{y})$
- K2:**  $\forall x : \varphi(x) \rightarrow \exists^{\leq k} \bar{y} : \psi(x, \bar{y})$
- K3:**  $\forall x : \varphi(x) \rightarrow \exists^{\geq k} \bar{y} : \psi(x, \bar{y})$

where  $x$  is a variable that ranges over nodes and values,  $\varphi(x)$  and  $\psi(x, \bar{y})$  are formulas of the form  $\exists \bar{z} : \xi$  with  $\bar{z}$  a vector of node and value variables and  $\xi$  a conjunction of atoms of the form  $\pi(z_i, z_j)$  with  $z_i$  and  $z_j$  either equal to  $x$ , in  $\bar{y}$ , or in  $\bar{z}$ , and  $\pi$  a standard PG-path expression, such that the free variables in  $\varphi(x)$  are just  $x$  and those in  $\psi(x, \bar{y})$  are  $x$  and the variables in  $\bar{y}$ .

The intuition of the constraints of the form **K1** is that they define a key constraint for all values or nodes selected by  $\varphi(x)$ . It states for such  $x$  that any vector of values and/or nodes  $\bar{y}$  that satisfies  $\psi(x, \bar{y})$  identifies  $x$ , i.e., is associated with at most one such  $x$ . So the symbol  $\Leftarrow$  should be read here as stating that the left-hand side is functionally determined by the right-hand side. More formally, its semantics is defined as being equivalent with the formula  $\forall \bar{y} : \exists^{\leq 1} x : \varphi(x) \wedge \psi(x, \bar{y})$ .

For the constraints of the forms **K2** and **K3** the interpretation is simply the usual one in first-order logic.

## D.5 The relationship with PG-Schema on full Property Graphs

The presented definitions introduce two important simplification w.r.t./ the original in [2]: (1) It is defined over common graphs which simplifies the property graph data model in several ways and (2) it assumes what is called the STRICT semantics of a graph type in [2] and ignores the LOOSE semantics. We briefly discuss here why these simplification preserve the essential characteristics of the original schema language.

**D.5.1 Concerning the simplification of the data model.** Recall that common graphs simplify property graphs in three ways: (1) nodes only have properties and no labels, (2) edges only have one label and no properties, and (3) edges have no independent identity. However, these features can be readily simulated in the common graph data model. For example, edges with identity can be simulated by nodes that have an outgoing edge with label *source* and an outgoing edge with label *target* to respectively the source node and the target node of the simulated edge. Moreover, labels can be simulated by introducing a special dummy value  $\Lambda$  that is used for keys that represent labels. For example, a node  $n$  where  $\rho(n)$  contains the pairs  $(Person, \Lambda)$ ,  $(Employee, \Lambda)$ ,  $(hiringDate, 12-Dec-2023)$ , and  $(fulltime, true)$ , simulates a node with labels *Person* and *Employee*, and properties *hiringDate* and *fulltime*.

It is not hard to see how under such a simulation the PG-Schema presented in this section could simulate a more powerful schema language where we could use tests in path expressions for the presence (or absence) of (combinations of) labels in path expressions and tests for presence (or absence) of (combinations of) properties of edges. Moreover, we could navigating via a simulated an edge and test for certain properties with a path expression of the form



1857  $source^- \cdot \pi \cdot target$  where  $\pi$  simulates any test over the content  
 1858 of the edge. Finally, we could straightforwardly simulate key and  
 1859 cardinality constraints for edges.

1860  
 1861 **D.5.2 Concerning the STRICT and LOOSE semantics.** In PG-Schema  
 1862 as described in [2] there is a separate LOOSE semantics defined for  
 1863 graph types. In that case the set of nodes types and the set of edge  
 1864 types in the graph type are ignored and a property graph is said to  
 1865 be already valid wrt. a graph type if it satisfies all graph constraints  
 1866 in the graph type. Note that in the original PG-Schema this did  
 1867 not mean that the sets of nodes types and edge types are entirely  
 1868 ignored, since there they also have the function of defining short-  
 1869 hands that can be used in the definitions of the graph constraints.

1870 The LOOSE interpretation can be easily simulated in the version  
 1871 presented here by letting the set of node types consist of  $\top$ , the  
 1872 trivial node type, and the set of edge types consist of  $\top \xrightarrow{\star} \top$ , the  
 1873 trivial edge type.

## 1874 D.6 The relationship with PG-Schema in the 1875 body of the paper

1876 The constraints of the forms **K2** and **K3** are very similar to the  
 1877 selector-shape pairs presented for PG-Schema in Section 5. Indeed,  
 1878 the selector is in this case the formula  $\varphi(x)$  and the shape are the  
 1879 formulas of the forms  $\exists^{\leq k} \bar{y} : \psi(x, \bar{y})$  and  $\exists^{\geq k} \bar{y} : \psi(x, \bar{y})$ . However,  
 1880 there are also several notable differences: (1) The schema only  
 1881 consists of constraints and does not separately mention node or  
 1882 edge types. (2) There are no edge types in path expressions. (3)  
 1883 All constraints are restricted so that  $\bar{y}$  is just a single variable. (4)  
 1884 There are no constraints of the form **K1**. (5) The constraints are  
 1885 syntactically restricted such that  $\varphi(x)$  is restricted to just one atom,  
 1886 so the form  $\exists z : \pi(x, z)$ , and  $\psi(x, \bar{y})$  is restricted to just one atom,  
 1887 so the form  $\pi(x, y)$ . Apart from these restriction, there is also a  
 1888 generalisation, since in Section 5 the shape expression are closed  
 1889 under intersection. That this does not change the expressive power  
 1890 is easy to see, since a selector-shape pair of the form  $(sel, (\varphi_1 \wedge \varphi_2))$   
 1891 can always be replaced with the combination of the pairs  $(sel, \varphi_1)$   
 1892 and  $(sel, \varphi_2)$  without changing the semantics of the schema.

1893 In the following subsections we discuss the previously mentioned  
 1894 restrictions.

1895  
 1896  
 1897 **D.6.1 No separate sets of node types and edge types.** It is not hard  
 1898 to show that this can be simulated. Assume for example we have a  
 1899 graph type with a set of node types  $\{\mathfrak{n}_1, \mathfrak{n}_2, \mathfrak{n}_3\}$ . The check that  
 1900 each node must be in the semantics of at least one of these node  
 1901 types can be simulated by the constraint

$$1902 \quad \forall x : \top(x, x) \rightarrow \exists y : (\mathfrak{n}_1 \mid \mathfrak{n}_2 \mid \mathfrak{n}_3)(x, y)$$

1903 Since node types are closed under the  $\mid$  operator,  $(\mathfrak{n}_1 \mid \mathfrak{n}_2 \mid \mathfrak{n}_3)$   
 1904 is a node type, and therefore an allowed path expression. Recall  
 1905 that a node type acts in a path expression as the identity relation  
 1906 restricted to nodes that are in the semantics of that type.

1907 Similarly, if the set of edge types of a graph type is  $\{\mathfrak{e}_1, \mathfrak{e}_2, \mathfrak{e}_3\}$ ,  
 1908 we can ensure that each edge is in the semantics of at least one of  
 1909 these edge types using the constraint

$$1910 \quad \forall x : \top(x, x) \rightarrow \exists^{\leq 0} y : \neg(\mathfrak{e}_1 \mid \mathfrak{e}_2 \mid \mathfrak{e}_3)(x, y).$$

1915 **D.6.2 No edge types in path expressions.** It is not hard to show that  
 1916 path expressions that contains tests involving edge types can be  
 1917 rewritten to equivalent path expressions that do not use edge types.

1918 We first consider the non-negated edge types in path expressions.  
 1919 We start with the observation that we can normalise edge types to  
 1920 a union of edge types that do not contain the  $\mid$  operator. This is  
 1921 based on the following equivalences for path semantics that allow  
 1922 us to push down the  $\mid$  operator:

- 1923 •  $(\mathfrak{n}_1 \mid \mathfrak{n}_2) \xrightarrow{\alpha} \mathfrak{n}_3 \equiv (\mathfrak{n}_1 \xrightarrow{\alpha} \mathfrak{n}_3 \mid \mathfrak{n}_2 \xrightarrow{\alpha} \mathfrak{n}_3)$
- 1924 •  $\mathfrak{n}_1 \xrightarrow{\alpha} (\mathfrak{n}_2 \mid \mathfrak{n}_3) \equiv (\mathfrak{n}_1 \xrightarrow{\alpha} \mathfrak{n}_2 \mid \mathfrak{n}_1 \xrightarrow{\alpha} \mathfrak{n}_3)$

1925 In a next normalisation step we can remove bottom-up the  $\&$ -  
 1926 operator using the following rules, where we use the symbol  $\mathfrak{e}_\emptyset$  to  
 1927 denote the empty edge type:

- 1928 •  $(\mathfrak{n}_1 \xrightarrow{\alpha} \mathfrak{n}_2) \& (\mathfrak{n}_3 \xrightarrow{\beta} \mathfrak{n}_4) \equiv (\mathfrak{n}_1 \& \mathfrak{n}_3) \xrightarrow{\alpha \cap \beta} (\mathfrak{n}_2 \& \mathfrak{n}_4)$

1929 where  $\cap$  is defined such that (1)  $\star \cap P = P \cap \star = P$  for  $P \subseteq \mathcal{P}$ , and  
 1930 (2)  $P \cap Q = P \cap Q$  for  $P, Q \subseteq \mathcal{P}$ .

1931 As a final normalisation step we get rid of edge types  $\mathfrak{n}_1 \xrightarrow{P} \mathfrak{n}_2$   
 1932 where  $P$  contains two or more predicates, by applying the rule:

- 1933 •  $\mathfrak{n}_1 \xrightarrow{\{p_1, \dots, p_k\}} \mathfrak{n}_2 \equiv (\mathfrak{n}_1 \xrightarrow{\{p_1\}} \mathfrak{n}_2 \mid \dots \mid \mathfrak{n}_1 \xrightarrow{\{p_k\}} \mathfrak{n}_2)$

1934 After these normalisation steps we will have rewritten the edge  
 1935 type to the form  $(\mathfrak{e}_1 \mid \dots \mid \mathfrak{e}_k)$  with each  $etype_i$  of one of the  
 1936 following forms: (1)  $\mathfrak{n}_1 \xrightarrow{\star} \mathfrak{n}_2$ , (2)  $\mathfrak{n}_1 \xrightarrow{\{p\}} \mathfrak{n}_2$ , and (3)  $\mathfrak{n}_1 \xrightarrow{\emptyset} \mathfrak{n}_2$ . We  
 1937 can express such an edge type  $(\mathfrak{e}_1 \mid \dots \mid \mathfrak{e}_k)$  as a path expression  
 1938  $(\pi_1 \cup \dots \cup \pi_k)$ , where each  $\pi_i$  is constructed as follows:

- 1939 •  $\mathfrak{n}_1 \xrightarrow{\star} \mathfrak{n}_2 \equiv \mathfrak{n}_1 \cdot \neg \emptyset \cdot \mathfrak{n}_2$
- 1940 •  $\mathfrak{n}_1 \xrightarrow{\{p\}} \mathfrak{n}_2 \equiv \mathfrak{n}_1 \cdot p \cdot \mathfrak{n}_2$
- 1941 •  $\mathfrak{n}_1 \xrightarrow{\emptyset} \mathfrak{n}_2 \equiv \neg \top$

1942 Recall that  $\neg \top$  is the negation of the trivial node type and so in a  
 1943 path expression represents the empty binary relation.

1944 We now turn our attention to negated edge types. These are  
 1945 normalised in the same as we described before, and we end up  
 1946 with an edge type of the form  $\neg(\mathfrak{e}_1 \mid \dots \mid \mathfrak{e}_k)$  with each  $\mathfrak{e}_i$  a  
 1947 primitive edge type. This can also be represented as a union of path  
 1948 expressions  $(\pi_1 \cup \dots \cup \pi_m)$  where each  $\pi_j$  is a path expression the  
 1949 expresses one way that the edge does not conform to any of the  
 1950 types in  $\mathfrak{e}_1, \dots, \mathfrak{e}_k$ . To illustrate this consider as an example the  
 1951 following negated edge type:

$$1952 \quad \neg(\mathfrak{n}_1 \xrightarrow{\{p\}} \mathfrak{n}_2 \mid \mathfrak{n}'_1 \xrightarrow{\{p'\}} \mathfrak{n}'_2)$$

1953 This can be simulated in a path expression by replacing it with the  
 1954 following path expression:

$$1955 \quad \neg \mathfrak{n}_1 \cdot \neg \mathfrak{n}'_1 \cdot \neg \emptyset \cup \neg \mathfrak{n}_1 \cdot \neg \{p'\} \cup \neg \mathfrak{n}_1 \cdot \neg \emptyset \cdot \neg \mathfrak{n}'_2 \cup$$

$$1956 \quad \cup \neg \mathfrak{n}'_1 \cdot \neg \{p\} \cup \neg \{p, p'\} \cup \neg \{p\} \cdot \neg \mathfrak{n}'_2 \cup$$

$$1957 \quad \cup \neg \mathfrak{n}'_1 \cdot \neg \emptyset \cdot \neg \mathfrak{n}_2 \cup \neg \{p'\} \cdot \neg \mathfrak{n}_2 \cup \neg \emptyset \cdot \neg \mathfrak{n}_2 \cdot \neg \mathfrak{n}'_2$$

1958 Note that this indeed enumerates all the ways that an edge could  
 1959 not be in the semantics of  $(\mathfrak{n}_1 \xrightarrow{\{p\}} \mathfrak{n}_2 \mid \mathfrak{n}'_1 \xrightarrow{\{p'\}} \mathfrak{n}'_2)$ . Basically we  
 1960 pick for each of the primitive edge types whether the edge is not  
 1961 in the semantics because of (1) the source node, (2) the label, or (3)  
 1962 the target node.

**D.6.3 Only single variable counting.** The restriction to allow only one variable in  $\bar{y}$  is introduced because in SHACL and ShEx all the counting is also restricted to single values and nodes, rather than tuples of values and nodes. Although this is often useful in real-world data modelling, e.g., to represent composite keys, this restriction is introduced to make PG-Schema more comparable to SHACL and ShEx.

**D.6.4 No  $\mathbf{K1}$  constraints.** After restricting to single-variable counting,  $\mathbf{K1}$  constraints are of the form

$$\forall x : \varphi(x) \Leftarrow \mathbf{Key} y : \psi(x, y).$$

Their semantics is defined by the formula  $\forall \bar{y} : \exists^{\leq 1} x : \varphi(x) \wedge \psi(x, \bar{y})$ .

If  $y$  matches nodes (which can be detected based on path expressions used in the atoms involving  $y$ ), we can equivalently express this constraint as

$$\forall y : \top(y, y) \rightarrow \exists^{\leq 1} x : \varphi(x) \wedge \psi(x, y)$$

If  $y$  matches values, then it is used in the first position of an atom whose path expressions begins from  $k^-$  or in the second position of an atom whose path expression ends with  $k$ . In either case, we can equivalently express this constraint as

$$\forall y : (k^- \cdot k)(y', y) \rightarrow \exists^{\leq 1} x : \varphi(x) \wedge \psi(x, y)$$

**D.6.5 Only one atom in formulas.** This restriction of the query language underlying PG-Schema limits the expressive power of PG-Schema, but similar restrictions would be imposed anyway on the Common Graph Schema Language by the limitations of SHACL and ShEx. Some additional expressive power could be gained by allowing tree-shaped conjunctions of atoms with at most 2 free variables, but this would further complicate the formal development.

## E MORE ON THE CORE

In this section we prove Proposition 1. Recall that common shapes are defined by the grammar

$$\begin{aligned} \varphi &::= \exists \pi \mid \exists^{\leq n} \pi_1 \mid \exists^{\geq n} \pi_1 \mid \exists c \wedge \nexists \neg P \mid \varphi \wedge \varphi . \\ c &::= \{ \} \mid \{ k : w \} \mid c \& c \mid c \mid c . \\ \pi_0 &::= \{ k : c \} \mid \neg \{ k : c \} \mid c \& \top \mid \neg(c \& \top) \mid \pi_0 \cdot \pi_0 . \\ \pi_1 &::= \pi_0 \cdot p \cdot \pi_0 \mid \pi_0 \cdot p^- \cdot \pi_0 \mid \pi_0 \cdot k \mid k^- \cdot \pi_0 . \\ \bar{\pi} &::= \pi_0 \mid p \mid \bar{\pi}^- \mid \bar{\pi} \cdot \bar{\pi} \mid \bar{\pi} \cup \bar{\pi} . \\ \pi &::= \bar{\pi} \mid \bar{\pi} \cdot k \mid k^- \cdot \bar{\pi} \mid k^- \cdot \bar{\pi} \cdot k' . \end{aligned}$$

where  $n \in \mathbb{N}$ ,  $P \subseteq_{fin} \mathcal{P}$ ,  $k, k' \in \mathcal{K}$ ,  $c \in \mathcal{V}$ , and  $p \in \mathcal{P}$ . We will refer to PG-path expressions defined by the nonterminal  $\pi_0$  in the grammar as *filters*.

The following two subsections describe the translations of common schemas to SHACL and ShEx. The translations are very similar but we include them both for the convenience of the reader.

### E.1 Translation to SHACL

**Lemma 1.** *For each open content type  $c$  there is a SHACL shape  $\varphi_c$  such that  $\mathcal{G}, v \models \varphi_c$  iff  $\rho(v) \in \llbracket c \rrbracket$  for all  $\mathcal{G} = (E, \rho)$  and  $v \in \text{Nodes}(\mathcal{G})$ .*

**PROOF.** For the content type  $\top$  the corresponding SHACL shape is  $\top$ . For a content type of the form

$$\{k_1 : w_1\} \& \{k_2 : w_2\} \& \dots \& \{k_m : w_m\} \& \top,$$

the corresponding SHACL shape is

$$\exists k_1.\text{test}(w_1) \wedge \exists k_2.\text{test}(w_2) \wedge \dots \wedge \exists k_m.\text{test}(w_m).$$

Finally, every open content type different from  $\top$  can be expressed as

$$(c_1 \mid \dots \mid c_\ell) \& \top,$$

where each  $c_i$  is a content type of the form  $\{k_1 : w_1\} \& \{k_2 : w_2\} \& \dots \& \{k_m : w_m\}$  for some  $m$ . The corresponding SHACL shape is

$$\varphi_1 \vee \dots \vee \varphi_\ell,$$

where  $\varphi_i$  is the SHACL shape corresponding to the content type  $c_i \& \top$ .  $\square$

**Lemma 2.** *For each filter  $\pi_0$  there is a SHACL shape  $\varphi_{\pi_0}$  such that  $\mathcal{G}, v \models \varphi_{\pi_0}$  iff  $(v, v) \in \llbracket \pi_0 \rrbracket^{\mathcal{G}}$  for all  $\mathcal{G}$  and  $v \in \text{Nodes}(\mathcal{G})$ .*

**PROOF.** By Lemma 1, the claim holds for  $\pi_0 = c \& \top$ . For  $\{k : c\}$  the corresponding SHACL shape is  $\exists k.\text{is}(w)$ . As SHACL shapes are closed under negation, the claim holds for  $\neg\{k : c\}$  and  $\neg(c \& \top)$ . Finally, concatenations of filters correspond to conjunctions of shapes, so the claim follows because SHACL shapes are closed under conjunction.  $\square$

**Lemma 3.** *For each common shape of the form  $\exists \pi$  there is a SHACL shape  $\varphi_{\exists \pi}$  such that  $\mathcal{G}, v \models \varphi_{\exists \pi}$  iff  $\mathcal{G}, v \models \exists \pi$  for all  $\mathcal{G}$  and  $v \in \text{Nodes}(\mathcal{G}) \cup \text{Values}(\mathcal{G})$ .*

**PROOF.** Let us first look at common shapes of the form  $\exists \pi$  where  $\pi$  is a concatenation of filters and atomic path expressions of the form  $p$ ,  $p^-$ ,  $k$ , or  $k^-$ . Without loss of generality we can assume that the concatenation ends with a filter or with  $k$ . We proceed by induction on the length of the concatenation. The base cases are  $\exists \pi_0$  and  $\exists k$ , which correspond to  $\varphi_{\pi_0}$  (Lemma 2) and  $\exists k.\top$ . For  $\exists \pi_0 \cdot \pi$  we can take  $\varphi_{\pi_0} \wedge \varphi_{\exists \pi}$ . For  $\exists p \cdot \pi$  we can take  $\exists p.\varphi_{\exists \pi}$ , and similarly for  $\exists p^- \cdot \pi$  and  $\exists k^- \cdot \pi$ .

The general case follows because SHACL shapes are closed under union. Indeed, because our PG-path expressions are star-free, we can assume without loss of generality that in each common shape of the form  $\exists \pi$ , the PG-path expression  $\bar{\pi}$  underlying  $\pi$  is a union of concatenations of filters and atomic path expressions of the form  $p$  or  $p^-$ . Then, for

$$\exists k^- \cdot (\pi^1 \cup \dots \cup \pi^m) \cdot k'$$

we can take

$$\varphi_{\exists k^- \cdot \pi^1 \cdot k'} \vee \dots \vee \varphi_{\exists k^- \cdot \pi^m \cdot k'}$$

Similarly for  $\exists k^- \cdot (\pi^1 \cup \dots \cup \pi^m)$ ,  $\exists (\pi^1 \cup \dots \cup \pi^m) \cdot k'$ , and  $\exists (\pi^1 \cup \dots \cup \pi^m)$ .  $\square$

**Lemma 4.** *For each common shape  $\varphi$  there is a SHACL shape  $\hat{\varphi}$  such that  $\mathcal{G}, v \models \varphi$  iff  $\mathcal{G}, v \models \hat{\varphi}$  for all  $\mathcal{G}$  and  $v \in \text{Nodes}(\mathcal{G}) \cup \text{Values}(\mathcal{G})$ .*

**PROOF.** Because SHACL shapes are closed under conjunction, it suffices to prove the claim for the atomic common shapes of the forms  $\exists \pi$ ,  $\exists^{\leq n} \pi_1$ ,  $\exists^{\geq n} \pi_1$ , and  $\exists c \wedge \nexists \neg P$ . The first case follows from Lemma 3.

Let us look at common shapes of the form  $\exists^{\geq n} \pi_1$ . If  $n = 0$  we can simply take  $\top$ . Suppose  $n > 0$ . Then, for

$$\exists^{\geq n} \pi_0 \cdot p \cdot \pi'_0$$

we can take

$$\varphi_{\pi_0} \wedge \exists^{\geq n} p \cdot \varphi_{\pi'_0},$$

and similarly for  $\exists^{\geq n} \pi_0 \cdot p^- \cdot \pi'_0$ ,  $\exists^{\geq n} \pi_0 \cdot k^- \cdot \pi'_0$ , and  $\exists^{\geq n} \pi_0 \cdot k$  (using  $\top$  instead of  $\varphi_{\pi'_0}$ ).

Next, we consider common shapes of the form  $\exists^{\leq n} \pi_1$ . For

$$\exists^{\leq n} \pi_0 \cdot p \cdot \pi'_0$$

we can take

$$\neg \varphi_{\pi_0} \vee \exists^{\leq n} p \cdot \varphi_{\pi'_0},$$

and similarly for  $\exists^{\leq n} \pi_0 \cdot p^- \cdot \pi'_0$ ,  $\exists^{\leq n} \pi_0 \cdot k^- \cdot \pi'_0$ , and  $\exists^{\leq n} \pi_0 \cdot k$  (again, using  $\top$  instead of  $\varphi_{\pi'_0}$ ).

Finally, let us consider a common shape of the form  $\exists \mathbb{C} \wedge \nexists \neg P$ . Suppose first that

$$\mathbb{C} = \{ \}.$$

Then, the corresponding SHACL shape is simply

$$\text{closed}(P).$$

Next, suppose that

$$\mathbb{C} = \{k_1 : \mathbb{V}_1\} \& \dots \& \{k_m : \mathbb{V}_m\}.$$

Then, the corresponding SHACL shape is

$$\exists k_1.\text{test}(\mathbb{V}_1) \wedge \dots \wedge \exists k_m.\text{test}(\mathbb{V}_m) \wedge \text{closed}(\{k_1, \dots, k_m\} \cup P).$$

In general, as in Lemma 1, we can assume that

$$\mathbb{C} = \mathbb{C}_1 \mid \dots \mid \mathbb{C}_m$$

where each  $\mathbb{C}_i$  is of one of the two forms considered above. The corresponding SHACL shape is then

$$\varphi_1 \vee \dots \vee \varphi_m$$

where  $\varphi_i$  is the SHACL shape corresponding to  $\exists \mathbb{C}_i \wedge \nexists \neg P$ , obtained as described above.  $\square$

**Lemma 5.** *For every common schema there is an equivalent SHACL schema.*

**PROOF.** Let  $\mathcal{S}$  be a common schema. We obtain an equivalent SHACL schema  $\mathcal{S}'$  by translating each  $(sel, \varphi) \in \mathcal{S}$  to  $(sel', \varphi')$  such that for all  $\mathcal{G}$  and  $v \in \mathcal{N} \cup \mathcal{V}$ ,

$$\begin{aligned} \mathcal{G}, v \models sel \text{ implies } \mathcal{G}, v \models \varphi \\ \text{iff} \\ \mathcal{G}, v \models sel' \text{ implies } \mathcal{G}, v \models \varphi'. \end{aligned}$$

Recall that  $sel$  is a common shape of one of the following forms:

$$\exists k, \exists p \cdot \pi, \exists p^- \cdot \pi, \exists \{k : c\} \cdot \pi, \exists (\{k : \mathbb{V}\} \& \top) \cdot \pi, \exists k^- \cdot \pi.$$

For  $sel'$  we take, respectively,

$$\exists k.\top, \exists p.\top, \exists p^-\top, \exists k.\top, \exists k.\top, \exists k^-\top.$$

For  $\varphi'$  we take  $\neg \varphi_{sel} \vee \hat{\varphi}$  where  $\varphi_{sel}$  is obtained using Lemma 3 and  $\hat{\varphi}$  is obtained using Lemma 4.  $\square$

## E.2 Translation to ShEx

**Lemma 6.** *For each open content type  $\mathbb{C}$  there is a ShEx shape  $\varphi_{\mathbb{C}}$  such that  $\mathcal{G}, v \models \varphi_{\mathbb{C}}$  iff  $\rho(v) \in \llbracket \mathbb{C} \rrbracket$  for all  $\mathcal{G} = (E, \rho)$  and  $v \in \text{Nodes}(\mathcal{G})$ .*

**PROOF.** For the content type  $\top$  the corresponding ShEx shape is  $\{\top\}$ .

For a content type of the form

$$\{k_1 : \mathbb{V}_1\} \& \dots \& \{k_m : \mathbb{V}_m\} \& \top,$$

the corresponding ShEx shape is

$$\{k_1.\text{test}(\mathbb{V}_1); \top\} \wedge \dots \wedge \{k_m.\text{test}(\mathbb{V}_m); \top\}.$$

Finally, every other open content type can be expressed as

$$(\mathbb{C}_1 \mid \dots \mid \mathbb{C}_\ell) \& \top,$$

where each  $\mathbb{C}_i$  has the form  $\{k_1 : \mathbb{V}_1\} \& \dots \& \{k_m : \mathbb{V}_m\}$  for some  $m$ . The corresponding ShEx shape is

$$\varphi_1 \vee \dots \vee \varphi_\ell,$$

where  $\varphi_i$  is the ShEx shape corresponding to the content type  $\mathbb{C}_i \& \top$ .  $\square$

**Lemma 7.** *For each filter  $\pi_0$  there is a ShEx shape  $\varphi_{\pi_0}$  such that  $\mathcal{G}, v \models \varphi_{\pi_0}$  iff  $(v, v) \in \llbracket \pi_0 \rrbracket^{\mathcal{G}}$  for all  $\mathcal{G}$  and  $v \in \text{Nodes}(\mathcal{G})$ .*

**PROOF.** By Lemma 6, the claim holds for  $\pi_0 = \mathbb{C} \& \top$ . For  $\{k : c\}$  the corresponding ShEx shape is  $\{k.\text{is}(w)\}$ . Because ShEx shapes are closed under negation, the claim also holds for  $\neg\{k : c\}$  and  $\neg(\mathbb{C} \& \top)$ . Finally, concatenations of filters correspond to conjunctions of shapes, so the claim follows because ShEx shapes are closed under conjunction.  $\square$

**Lemma 8.** *For each common shape of the form  $\exists \pi$  there is a ShEx shape  $\varphi_{\exists \pi}$  such that  $\mathcal{G}, v \models \varphi_{\exists \pi}$  iff  $\mathcal{G}, v \models \exists \pi$  for all  $\mathcal{G}$  and  $v \in \text{Nodes}(\mathcal{G}) \cup \text{Values}(\mathcal{G})$ .*

**PROOF.** Let us first look at common shapes of the form  $\exists \pi$  where  $\pi$  is a concatenation of filters and atomic path expressions of the form  $p, p^-, k$ , or  $k^-$ . Without loss of generality we can assume that the concatenation ends with a filter or with  $k$ . We proceed by induction on the length of the concatenation. The base cases are  $\exists \pi_0$  and  $\exists k$ , which correspond to  $\varphi_{\pi_0}$  (Lemma 7) and  $\{k.\top\}; \top$ , respectively. For  $\exists \pi_0 \cdot \pi$  we can take  $\varphi_{\pi_0} \wedge \varphi_{\exists \pi}$ . For  $\exists p \cdot \pi$  we can take  $\{p.\varphi_{\exists \pi}; \top\}$ , and similarly for  $\exists p^- \cdot \pi$  and  $\exists k^- \cdot \pi$ .

The general case follows because ShEx shapes are closed under union. Indeed, because our PG-path expressions are star-free, we can assume without loss of generality that in each common shape of the form  $\exists \pi$ , the PG-path expression  $\bar{\pi}$  underlying  $\pi$  is a union of concatenations of filters and atomic path expressions of the form  $p$  or  $p^-$ . Then, for

$$\exists k^- \cdot (\pi^1 \cup \dots \cup \pi^m) \cdot k'$$

we can take

$$\varphi_{\exists k^- \cdot \pi^1 \cdot k'} \vee \dots \vee \varphi_{\exists k^- \cdot \pi^m \cdot k'}.$$

Similarly for  $\exists k^- \cdot (\pi^1 \cup \dots \cup \pi^m)$ ,  $\exists (\pi^1 \cup \dots \cup \pi^m) \cdot k'$ , and  $\exists (\pi^1 \cup \dots \cup \pi^m)$ .  $\square$

**Lemma 9.** *For each common shape  $\varphi$  there is a ShEx shape  $\hat{\varphi}$  such that  $\mathcal{G}, v \models \varphi$  iff  $\mathcal{G}, v \models \hat{\varphi}$  for all  $\mathcal{G}$  and  $v \in \text{Nodes}(\mathcal{G}) \cup \text{Values}(\mathcal{G})$ .*

PROOF. Because ShEx shapes are closed under conjunction, it suffices to prove the claim for the atomic common shapes of the forms  $\exists \pi$ ,  $\exists^{\leq n} \pi_1$ ,  $\exists^{\geq n} \pi_1$ , and  $\exists \mathbb{C} \wedge \nexists \neg P$ . The first case follows from Lemma 8.

Let us look at common shapes of the form  $\exists^{\geq n} \pi_1$ . If  $n = 0$  we can simply take  $\{\top\}$ . Suppose  $n > 0$ . Then, for

$$\exists^{\geq n} \pi_0 \cdot p \cdot \pi'_0$$

we can take

$$\varphi_{\pi_0} \wedge \{(p \cdot \varphi_{\pi'_0})^n; \top\},$$

and similarly for  $\exists^{\geq n} \pi_0 \cdot p^- \cdot \pi'_0$ ,  $\exists^{\geq n} \pi_0 \cdot k^- \cdot \pi'_0$ , and  $\exists^{\geq n} \pi_0 \cdot k$  (using  $\{\top\}$  instead of  $\varphi_{\pi'_0}$ ).

Next, we consider common shapes of the form  $\exists^{\leq n} \pi_1$ . For

$$\exists^{\leq n} \pi_0 \cdot p \cdot \pi'_0$$

we can take

$$\neg \varphi_{\pi_0} \vee \neg \{(p \cdot \varphi_{\pi'_0})^{n+1}; \top\}$$

and similarly for  $\exists^{\leq n} \pi_0 \cdot p^- \cdot \pi'_0$ ,  $\exists^{\leq n} \pi_0 \cdot k^- \cdot \pi'_0$ , and  $\exists^{\leq n} \pi_0 \cdot k$  (again, using  $\{\top\}$  instead of  $\varphi_{\pi'_0}$ ).

Before we move on, let us introduce a bit of syntactic sugar. For a set  $Q = \{q_1, q_2, \dots, q_n\} \subseteq \mathcal{P} \cup \mathcal{K}$  we write  $Q^*$  for the triple expression  $(q_1 \cdot \{\top\} \mid q_2 \cdot \{\top\} \mid \dots \mid q_n \cdot \{\top\})^*$ .

We are now ready to consider a common shape of the form  $\exists \mathbb{C} \wedge \nexists \neg P$ . Suppose first that

$$\mathbb{C} = \{\}.$$

Then, the corresponding ShEx shape is simply

$$\{P^*; (\neg \emptyset^-)^*\}.$$

Next, suppose that

$$\mathbb{C} = \{k_1 : \mathbb{V}_1\} \& \dots \& \{k_m : \mathbb{V}_m\}.$$

Then, the corresponding ShEx shape is

$$\varphi_{\mathbb{C} \& \top} \wedge \{ \{k_1, \dots, k_m\}^*; P^*; (\neg \emptyset^-)^* \},$$

where  $\varphi_{\mathbb{C} \& \top}$  is obtained from Lemma 6. In general, as in Lemma 6, we can assume that

$$\mathbb{C} = \mathbb{C}_1 \mid \dots \mid \mathbb{C}_m$$

where each  $\mathbb{C}_i$  is of one of the two forms considered above. The corresponding ShEx shape is then

$$\varphi_1 \vee \dots \vee \varphi_m$$

where  $\varphi_i$  is the ShEx shape corresponding to  $\exists \mathbb{C}_i \wedge \nexists \neg P$ , obtained as described above.  $\square$

**Lemma 10.** *For every common schema there is an equivalent ShEx schema.*

PROOF. Let  $\mathcal{S}$  be a common schema. We obtain an equivalent ShEx schema  $\mathcal{S}'$  by translating each  $(sel, \varphi) \in \mathcal{S}$  to  $(sel', \varphi')$  such that for all  $\mathcal{G}$  and  $v \in \mathcal{N} \cup \mathcal{V}$ ,

$$\begin{aligned} \mathcal{G}, v \models sel \text{ implies } \mathcal{G}, v \models \varphi \\ \text{iff} \\ \mathcal{G}, v \models sel' \text{ implies } \mathcal{G}, v \models \varphi'. \end{aligned}$$

Recall that  $sel$  is a common shape of one of the following forms:

$$\exists k, \exists p \cdot \pi, \exists p^- \cdot \pi, \exists \{k : c\} \cdot \pi, \exists (\{k : \mathbb{V}\} \& \top) \cdot \pi, \exists k^- \cdot \pi.$$

If  $sel$  is of the form

$$\exists k, \exists \{k : c\} \cdot \pi, \text{ or } \exists (\{k : \mathbb{V}\} \& \top) \cdot \pi,$$

for  $sel'$  we take  $\{k \cdot \{\top\}; \top\}$ . In the remaining cases, we take, respectively,

$$\{p \cdot \{\top\}; \top\}, \{p^- \cdot \{\top\}; \top\}, \{k^- \cdot \{\top\}; \top\}.$$

For  $\varphi'$  we take  $\neg \varphi_{sel} \vee \hat{\varphi}$  where  $\varphi_{sel}$  is obtained using Lemma 8 and  $\hat{\varphi}$  is obtained using Lemma 9.  $\square$