# A Structural Subsumption based Similarity Measure for the Description Logic $\mathcal{ALEH}$

Boontawee Suntisrivaraporn and Suwan Tongphu

*School of Information, Computer and Communication Technology,*
*Sirindhorn International Institute of Technology, Thammasat University, Pathumthani, Thailand*

Abstract:      Description Logics (DLs) are a family of logic-based knowledge representation formalisms, which can be used to develop ontologies in a formally well-founded way. The standard reasoning service of subsumption has proved indispensable in ontology design and maintenance. This checks, relative to the logical definitions in the ontology, whether one concept is more general/specific than another. When no subsumption relationship is identified, however, no information about the two concepts can be given. This work extends from an existing work on similarity measure in $\mathcal{ELH}$ to the more expressive description logic $\mathcal{ALEH}$. We introduce generalizations of the notion of normalization and homomorphism in $\mathcal{ALEH}$ which are then employed at the heart of our semantic similarity measure. The proposed similarity measure computes a numerical degree of similarity between two $\mathcal{ALEH}$ concept descriptions despite not being in the subsumption relation.

## 1 INTRODUCTION

Representing knowledge base is one interesting topic in artificial intelligence. Among various techniques of semantic-level analysis, one commonly well-founded way is through the help of Description Logics (DLs) (Baader et al., 2007). Being recommended by W3C, DLs (i.e. the logical underpinning of the Web Ontology Language (OWL)) become a standard tool for formally and systematically modelling a knowledge base. Besides their unambiguous syntax and semantics which are essential for ontology modelling and sharing, DLs provide several useful reasoning services that allow inferencing of implicit knowledge from the one explicitly defined. For example, with a service of a subclass-superclass relation identification (concept subsumption), two defined concepts which are visually out of subsumption relation may be logically classified into the same hierarchy. Though seemingly useful, the classical DL reasoning service of concept subsumption merely produces a crisp response. The service indeed provides a positive conclusion if and only if all necessary and sufficient conditions of being in the subclass–superclass relation are fully satisfied. Otherwise, alas, it will suggest that the two concepts are irrelevant to each other.

In some concrete situation, checking for subsumption relation may not be adequate. Consider for exam-

ple the case in which a new disease, which is closely similar to the existing one, is being discovered. Since we know that the two diseases are similar, checking for their common characteristics would likely provide a beneficial clue to the disease etiology. Therefore, it would be easy to suggest an appropriate treatment from previously known diseases to another new one.

This work is an extension of an existing similarity measures for DLs in the $\mathcal{EL}$ family (Suntisrivaraporn, 2013; Tongphu and Suntisrivaraporn, 2014; Tongphu and Suntisrivaraporn, 2015) to the strictly more expressive DL $\mathcal{ALEH}$. The method is based on the known homomorphism-based structural subsumption and produces a numerical degree of similarity between two $\mathcal{ALEH}$ concept descriptions despite not being in the subsumption relation.

The rest of the paper is organized in order. The background on the DL $\mathcal{ALEH}$, unfoldable TBoxes, and the structural subsumption algorithm is presented in the next section. Section 3 and 4 introduce the notions of homomorphism degree and $\mathcal{ALEH}$ semantic similarity measure, respectively, and exemplify the introduced measure by means of a small yet prototypical medical ontology. Section 5 suggests a possible extension of similarity measure for the DL $\mathcal{ALCH}$. Related works are discussed in Section 6, and the last section gives some concluding remarks.

## 2 BACKGROUND

In DLs, *concept descriptions* are inductively defined with the help of a set of *constructors*, starting with a set CN of *concept names* and a set RN of *role names*. $\mathcal{ALEH}$ concept descriptions are formed using the constructors shown in the upper part of Table 1. An $\mathcal{ALEH}$ *terminology* or *TBox* is a finite set of concept definitions and role hierarchy axioms, of which the syntactic forms are shown in the lower part of Table 1. A TBox is called *unfoldable* if it is *definitorial* (i.e. containing at most one concept definition for each concept name), and *acyclic* (i.e. not containing cyclic dependencies). Figure 1 depicts an example unfoldable $\mathcal{ALEH}$ TBox. The set $\mathsf{CN}^{\mathsf{def}}$ of *defined concepts* comprises those concept names that appear on the left hand side of a concept definition. All other concept names are called *primitive concepts*, denoted by $\mathsf{CN}^{\mathsf{pri}}$. Since the DL $\mathcal{ALEH}$ allows for atomic negation, for convenience, we denote by $\mathsf{CN}^{\mathsf{label}}$ the set of all primitive concepts, their negations, and the bottom concept, i.e. $\mathsf{CN}^{\mathsf{label}} = \{A, \neg A \mid A \in \mathsf{CN}^{\mathsf{pri}}\} \cup \{\bot\}$. Conventionally, $r, s$ possibly with subscripts are used to range over RN, $A, B$ to range over CN, and $C, D$ to range over concept descriptions. Primitive concept definitions are commonly found in realistic terminologies to define those concepts, of which only necessary conditions are known. For instance,

$$\mathsf{HappyMan} \sqsubseteq \mathsf{Man} \sqcap \mathsf{Rich} \sqcap \exists \mathsf{child}.\mathsf{Beautiful} \tag{1}$$

Such a primitive definition $B \sqsubseteq D$ can easily be transformed into a semantically equivalent full definitions $B \equiv X \sqcap D$ where $X$ is a fresh concept name.

Like other DLs, the semantics of $\mathcal{ALEH}$ is defined in terms of *interpretations* $I = (\Delta^I, \cdot^I)$, where the domain $\Delta^I$ is a non-empty set of individuals, and the interpretation function $\cdot^I$ maps each concept name $A \in \mathsf{CN}$ to a subset $A^I$ of $\Delta^I$ and each role name $r \in \mathsf{RN}$ to a binary relation $r^I$ on $\Delta^I$. The extension of $\cdot^I$ to arbitrary concept descriptions is inductively defined, as shown in the semantics column of Table 1. An interpretation $I$ is a *model* of a TBox $O$ if, for each concept definition in $O$, the conditions given in the semantics column of Table 1 are satisfied. The main inference problem for $\mathcal{ALEH}$ is the subsumption problem:

**Definition 1** (Concept Subsumption)**.** *Given two* $\mathcal{ALEH}$ *concept descriptions* $C, D$ *and an* $\mathcal{ALEH}$ *TBox* $O$, $C$ *is subsumed by* $D$ *w.r.t.* $O$ *(written* $C \sqsubseteq_O D$*) if* $C^I \subseteq D^I$ *in every model* $I$ *of* $O$*. Moreover,* $C, D$ *are equivalent w.r.t.* $O$ *(written* $C \equiv_O D$*) if* $C \sqsubseteq_O D$ *and* $D \sqsubseteq_O C$*.*



| $\omega_1$ | Woman $\equiv$ Female $\sqcap$ Person |
| $\omega_2$ | Man $\equiv \neg$Female $\sqcap$ Person |
| $\omega_3$ | Parent $\equiv$ Person $\sqcap \exists$child.Person |
| $\omega_4$ | Mother $\equiv$ Woman $\sqcap$ Parent |
| $\omega_5$ | Father $\equiv$ Man $\sqcap$ Parent |
| $\omega_6$ | MotherNoSon $\equiv$ Mother $\sqcap \forall$child.Woman |
| $\omega_7$ | MotherNoDaughter $\equiv$ Mother $\sqcap \forall$child.Man |
| $\omega_8$ | FosterFather $\equiv$ Man $\sqcap \exists$fchild.Person |
| $\omega_9$ | NonFosterFather $\equiv$ Father $\sqcap \forall$fchild.$\bot$ |
| $\omega_{10}$ | fchild $\sqsubseteq$ child |

Figure 1: An example $\mathcal{ALEH}$ terminology $O_{\mathsf{family}}$; here child and fchild are shorthands for hasChild and hasFosterChild, respectively.

Provided that the TBox is unfoldable (i.e. acyclic and definitional), any $\mathcal{ALEH}$ concept description can be expanded to an equivalent one that may use any role names but consists only of primitive concept names, their negations and the bottom concept from $\mathsf{CN}^{\mathsf{label}}$. This can be done by repeatedly replacing a defined concept by its definition until no more defined concepts appear in the concept description. Consider, for instance, the concept MotherNoSon along with its definition $\omega_6$ in Figure 1. By replacing the defined concept Mother and Woman with their corresponding descriptions (see $\omega_4$ and $\omega_1$), the description can be expanded to:

$$\mathsf{Female} \sqcap \mathsf{Person} \sqcap \exists \mathsf{child}.\mathsf{Person} \sqcap \\ \forall \mathsf{child}.(\mathsf{Female} \sqcap \mathsf{Person}) \tag{2}$$

where $\mathsf{Person}, \mathsf{Female} \in \mathsf{CN}^{\mathsf{label}}$. We denote by $\hat{C}$ the expanded equivalence of the concept description $C$.

We can assume without loss of generality that an $\mathcal{ALEH}$ concept $C$ can be expanded and has the following form:

$$\prod_{i=1}^{l} L_i \quad \sqcap \quad \prod_{j=1}^{m} \exists r_j.D_j \quad \sqcap \quad \prod_{k=1}^{n} \forall s_k.E_k$$

where $L_i \in \mathsf{CN}^{\mathsf{label}}$, $r_j, s_k \in \mathsf{RN}$, and $D_j, E_k$ are $\mathcal{ALEH}$ concept descriptions in the same format as $C$. For simplicity, we assign $\mathcal{P}_C := \{L_1, \ldots, L_l\}$, $\mathcal{E}_C := \{\exists r_1.D_1, \ldots, \exists r_m.D_m\}$, and $\mathcal{A}_C := \{\forall s_1.E_1, \ldots, \forall s_n.E_n\}$. Also, we denote by $\mathcal{R}^{\exists r}$ and $\mathcal{R}^{\forall r}$ the sets of all super-roles and of all sub-roles of $r$, respectively. That is, $\mathcal{R}^{\exists r} = \{s \in \mathsf{RN} \mid r \sqsubseteq^* s\}$ and $\mathcal{R}^{\forall r} = \{t \in \mathsf{RN} \mid t \sqsubseteq^* r\}$ where where $\sqsubseteq^*$ represents the reflexive-transitive closure of $\sqsubseteq$ over role names. However, since a normalized $\mathcal{ALEH}$ concept description makes implicit description explicit and yet

Table 1: Syntax and semantics of the DL $\mathcal{ALEH}$ and DL $\mathcal{ALCH}$.

| Name | Syntax | Semantics | $\mathcal{ALEH}$ | $\mathcal{ALCH}$ |
|---|---|---|---|---|
| bottom | $\bot$ | $\emptyset$ | ✓ | ✓ |
| top | $\top$ | $\Delta^I$ | ✓ | ✓ |
| concept name | $A$ | $A^I \subseteq \Delta^I$ | ✓ | ✓ |
| atomic negation | $\neg A$ | $\Delta^I \backslash A$ | ✓ | ✓ |
| concept negation | $\neg C$ | $\Delta^I \backslash C$ | | ✓ |
| concept conjunction | $C \sqcap D$ | $C^I \cap D^I$ | ✓ | ✓ |
| concept disjunction | $C \sqcup D$ | $C^I \cup D^I$ | | ✓ |
| existential restriction | $\exists r.C$ | $\{x \mid \exists y \in \Delta^I : (x,y) \in r^I \wedge y \in C^I\}$ | ✓ | ✓ |
| value restriction | $\forall r.C$ | $\{x \mid \forall y \in \Delta^I : (x,y) \in r^I \Rightarrow y \in C^I\}$ | ✓ | ✓ |
| primitive definition | $B \sqsubseteq D$ | $A^I \subseteq D^I$ | ✓ | ✓ |
| full definition | $B \equiv D$ | $A^I = D^I$ | ✓ | ✓ |
| role hierarchy | $r \sqsubseteq s$ | $r^I \subseteq s^I$ | ✓ | ✓ |

preserves equivalence, we exhaustively apply the following normalization rules to the $\mathcal{ALEH}$ concept descriptions after expansion. The normalization rules below are modulo commutativity of conjunction:

$$
\begin{aligned}
\forall s.C \sqcap \forall r.D &\rightarrow \forall s.C \sqcap \forall r.(C \sqcap D) \\
\forall s.C \sqcap \exists r.D &\rightarrow \forall s.C \sqcap \exists r.(C \sqcap D) \\
\forall r.\top &\rightarrow \top \\
C \sqcap \top &\rightarrow C \\
A \sqcap \neg A &\rightarrow \bot \\
\exists r.\bot &\rightarrow \bot \\
C \sqcap \bot &\rightarrow \bot
\end{aligned}
$$

where $s \in \mathcal{R}^{\exists r}$. Note that the first two normalization rules generalize the corresponding ones in (Baader and Küsters, 2006) where a role hierarchy is taken into consideration. In fact, for a super-role $s$ of $r$, it is the case that $\forall s.C$ implies $\forall r.C$.

For example, let MotherNoSon be expanded and has the form as shown in Equation 2. By applying the above rules, a normalized concept description of MotherNoSon can be exemplified as follows:

Female $\sqcap$ Person $\sqcap$ $\exists$child.(Female $\sqcap$ Person) $\sqcap$ $\forall$child.(Female $\sqcap$ Person)

In (Baader and Küsters, 2000; Baader, 2003), a characterization of subsumption in $\mathcal{ALEH}$ w.r.t. an unfoldable TBox using homomorphism has been proposed. Instead of considering concept descriptions directly, the characterization considers so-called $\mathcal{ALEH}$ description trees that structurally correspond to the $\mathcal{ALEH}$ concept descriptions. Given the expanded concept description $C$, beginning from the top level, such a description can recursively be translated into an $\mathcal{ALEH}$ description tree $\mathcal{G}_C := (V, E, v_0, \ell, \rho)$

where $V$ is a set of nodes, $E \subseteq V \times V$ is a set of edges, $v_0 \in V$ is the root, $\ell : V \rightarrow 2^{\mathsf{CN^{label}}}$ is a node labelling function, and $\rho : E \rightarrow 2^{\mathsf{RN}}$ is an edge labelling function. The translation can be done using the following steps:

i. Assign $\mathcal{P}_C$ to $\ell(v_0)$.

ii. For each $\exists r.X_j \in \mathcal{E}_C$, introduce a new node $w$ to $V$, add an edge $(v_0, w)$ to $E$, and assign $\mathcal{R}^{\exists r}$ to $\rho(v_0, w)$. Repeat from step (i) by treating $w$ as $v_0$ and $X$ as $C$.

iii. For each $\forall r.Y_j \in \mathcal{A}_C$, introduce a new node $w'$ to $V$, add an edge $(v_0, w')$ to $E$, and assign $\mathcal{R}^{\forall r}$ to $\rho(v_0, w')$. Repeat from step (i) by treating $w'$ as $v_0$ and $Y$ as $C$.

In essence, the root $v_0$ of the $\mathcal{ALEH}$ description tree $\mathcal{G}_C$ has $\mathcal{P}_C$ as its label; has $m$ existential edges, each labeled with $R^{\exists r_j}$ to a vertex $w_j$; and has $n$ universal edges, each labeled with $R^{\forall s_k}$ to a vertex $w'_k$, for $1 \le j \le m$ and $1 \le k \le n$. Each of the child nodes $w_j$ and $w'_k$ is the root of a similar tree structure which forms a subtree of $\mathcal{G}_C$.

**Definition 2** (Homomorphism). *A homomorphism from an $\mathcal{ALEH}$ description tree $\mathcal{G} = (V, E, v_0, \ell, \rho)$ into an $\mathcal{ALEH}$ description tree $\mathcal{G}' = (V', E', v'_0, \ell', \rho')$ is a mapping $h : V \rightarrow V'$ such that:*

*i. $h(v_0) = v'_0$,*

*ii. $\ell(v) \subseteq \ell'(h(v))$ for all $v \in V$,*

*iii. for each existential edge $(v, w) \in E$ with $\rho(v, w) = \mathcal{R}^{\exists r}$, there exists $(h(v), h(w)) \in E'$ such that $\rho'(h(v), h(w)) = \mathcal{R}^{\exists s}$ and $\mathcal{R}^{\exists r} \subseteq \mathcal{R}^{\exists s}$, and*

*iv. for each universal edge $(v,w) \in E$ with $\rho(v,w) = \mathcal{R}^{\forall r}$, there exists $(h(v), h(w)) \in E'$ such that $\rho'(h(v), h(w)) = \mathcal{R}^{\forall t}$ and either of the following holds:*

*a. $\mathcal{R}^{\forall r} \subseteq \mathcal{R}^{\forall t}$, or*

*b. $h(v) = h(w)$ and $\ell'(h(v)) = \{\bot\}$.*

Observe that this generalizes the notion of homomorphism first introduced in (Baader and Küsters, 2006) by allowing each edge label to be a set of role names instead of a mere role name. Moreover, it simplifies the condition for existential edge mapping by omitting $\bot$ since any existential successor with $\bot$ as its label must be collapsed due to the normalization.

The subsumption is then characterized by means of an existence of a homomorphism in the reverse direction.

**Theorem 1** ((Baader and Küsters, 2006)). *Let $C, D$ be $\mathcal{ALEH}$ concept descriptions, and $\mathcal{G}_C, \mathcal{G}_D$ the corresponding $\mathcal{ALEH}$ concept description trees. Then, $C \sqsubseteq D$ iff there exists a homomorphism $h : \mathcal{G}_D \rightarrow \mathcal{G}_C$ which maps the root of $\mathcal{G}_D$ to the root of $\mathcal{G}_C$.*

Consider the normalized description for MotherNoSon as previously mentioned and the following normalized descriptions for Mother and NonFosterFather:

$$\text{Female} \sqcap \text{Person} \sqcap \exists \text{child.Person} \qquad (3)$$

$$\neg \text{Female} \sqcap \text{Person} \sqcap \exists \text{child.Person} \sqcap \forall \text{fchild.} \bot \qquad (4)$$

Figure 2 depicts the $\mathcal{ALEH}$ description trees $\mathcal{G}_{\text{NonFosterFather}}$ (left) $\mathcal{G}_{\text{Mother}}$ (center), and $\mathcal{G}_{\text{MotherNoSon}}$ (right). It is important to note here that $\mathcal{R}^{\forall \text{child}} = \{\text{fchild}, \text{child}\}$ and $\mathcal{R}^{\forall \text{fchild}} = \{\text{fchild}\}$ since $\omega_{10}$ is the only role hierarchy axiom in the ontology. This figure shows a homomorphism $h$ (dashed arrows) that maps the root $u_0$ of $\mathcal{G}_{\text{Mother}}$ to the root $v_0$ of $\mathcal{G}_{\text{MotherNoSon}}$. It also demonstrates a failed attempt to map (see the dotted arrow) the root of $\mathcal{G}_{\text{Mother}}$ to the root of $\mathcal{G}_{\text{NonFosterFather}}$. Theorem 1 ensures that MotherNoSon $\sqsubseteq_O$ Mother and NonFosterFather $\not\sqsubseteq_O$ Mother.

Though sharing some common features between MotherNoSon and NonFosterFather (i.e. both are Person ), the classical reasoning of subsumption cannot tell how similar the two descriptions are. This leads to an introduction of a concept similarity measure based on the structural characterization. Instead of merely giving either positive or negative result between descriptions, the proposed measure calculates a numerical value ranging between 0 and 1. Intuitively, the larger the number approaching to 1, the more similar the two concepts are.

# 3 HOMOMORPHISM DEGREE IN $\mathcal{ALEH}$

As suggested by Theorem 1, an existence of a homomorphism mapping from one $\mathcal{ALEH}$ description tree to another implies a subsumption relationship in a reverse direction. We extend the idea to the case where a homomorphism between the two $\mathcal{ALEH}$ description trees does not exist but there is a shared structure. Let $C, D$ be $\mathcal{ALEH}$ concept descriptions, and $\mathcal{G}_C$ and $\mathcal{G}_D$ be the corresponding $\mathcal{ALEH}$ description trees. Also, let $\mathcal{P}_C, \mathcal{P}_D, \mathcal{E}_C, \mathcal{E}_D, \mathcal{A}_C$, and $\mathcal{A}_D$ be as defined in the previous section. We define the homomorphism degree from $\mathcal{G}_D$ to $\mathcal{G}_C$ as follows:

**Definition 3** (Homomorphism Degree). *Let $\mathbf{G}^{\mathcal{ALEH}}$ be the set of all $\mathcal{ALEH}$ description trees. The homomorphism degree function $\text{hd} : \mathbf{G}^{\mathcal{ALEH}} \times \mathbf{G}^{\mathcal{ALEH}} \rightarrow [0,1]$ is defined as follows:*

$$\begin{aligned} \text{hd}(\mathcal{G}_D, \mathcal{G}_C) := \quad & (1 - \mu^e - \mu^a) \cdot \text{p-hd}(\mathcal{P}_D, \mathcal{P}_C) + \\ & \mu^e \cdot \text{e-set-hd}(\mathcal{E}_D, \mathcal{E}_C) + \\ & \mu^a \cdot \text{a-set-hd}(\mathcal{A}_D, \mathcal{A}_C) \end{aligned} \qquad (5)$$

*where $|\cdot|$ represents the set cardinality, $\mu^e = \frac{|\mathcal{E}_D|}{|\mathcal{P}_D \cup \mathcal{E}_D \cup \mathcal{A}_D|}$, and $\mu^a = \frac{|\mathcal{A}_D|}{|\mathcal{P}_D \cup \mathcal{E}_D \cup \mathcal{A}_D|}$;*

$$\text{p-hd}(\mathcal{P}_D, \mathcal{P}_C) := \begin{cases} 1 & \text{if } \mathcal{P}_D = \emptyset \text{ or } \mathcal{P}_C = \{\bot\} \\ \frac{|\mathcal{P}_D \cap \mathcal{P}_C|}{|\mathcal{P}_D|} & \text{otherwise,} \end{cases} \qquad (6)$$

$$\text{e-set-hd}(\mathcal{E}_D, \mathcal{E}_C) := \\ \begin{cases} 1 & \text{if } \mathcal{E}_D = \emptyset \\ 0 & \text{if } \mathcal{E}_D \neq \emptyset, \mathcal{E}_C = \emptyset \\ \sum_{\varepsilon_i \in \mathcal{E}_D} \frac{max\{\text{e-hd}(\varepsilon_i, \varepsilon_j) : \varepsilon_j \in \mathcal{E}_C\}}{|\mathcal{E}_D|} & \text{otherwise,} \end{cases} \qquad (7)$$

*where $\varepsilon_i, \varepsilon_j$ are existential restrictions;*

$$\text{e-hd}(\exists r.X, \exists s.Y) := \\ \gamma^e(\nu^e(r) + (1 - \nu^e(r)) \cdot \text{hd}(\mathcal{G}_X, \mathcal{G}_Y)) \qquad (8)$$

*where $\gamma^e = \frac{|\mathcal{R}^{\exists r} \cap \mathcal{R}^{\exists s}|}{|\mathcal{R}^{\exists r}|}$ and $\nu^e : \text{RN} \rightarrow [0,1)$.*

$$\text{a-set-hd}(\mathcal{A}_D, \mathcal{A}_C) := \\ \begin{cases} 1 & \text{if } \mathcal{A}_D = \emptyset, \\ 0 & \text{if } \mathcal{A}_D \neq \emptyset, \mathcal{A}_C = \emptyset, \\ \sum_{\alpha_i \in \mathcal{A}_D} \frac{max\{\text{a-hd}(\alpha_i, \alpha j) : \alpha j \in \mathcal{A}_C\}}{|\mathcal{A}_D|} & \text{otherwise} \end{cases} \qquad (9)$$

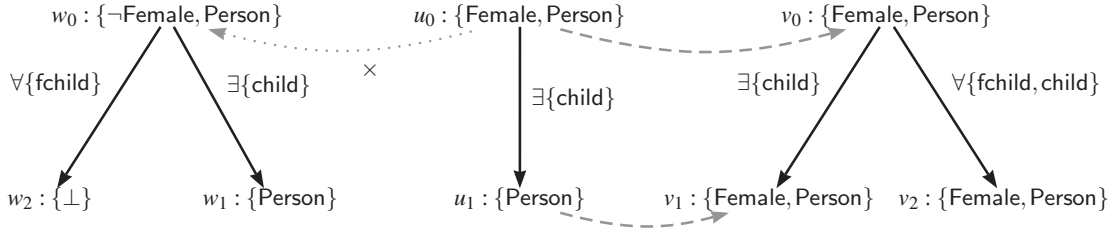*where $\alpha_i, \alpha_j$ are universal restrictions; and finally*

Figure 2: A homomorphism $h$ (dashed arrows) that maps the root of $\mathcal{G}_{\text{Mother}}$ to the root of $\mathcal{G}_{\text{MotherNoSon}}$; a failed attempt to identify a homomorphism (dotted arrows) that maps the root of $\mathcal{G}_{\text{Mother}}$ to the root of $\mathcal{G}_{\text{NonFosterFather}}$.

a-hd$(\forall r.X, \forall s.Y) :=$
$$\begin{cases} \gamma^{\mathsf{a}} & \text{if } \mathcal{P}_Y = \{\bot\}, \\ \gamma^{\mathsf{a}}(\nu^{\mathsf{a}}(r) + (1 - \nu^{\mathsf{a}}(r)) \cdot \text{hd}(\mathcal{G}_X, \mathcal{G}_Y)) & \text{otherwise} \end{cases}$$
(10)

where $\gamma^{\mathsf{a}} = \frac{|\mathcal{R}^{\forall r} \cap \mathcal{R}^{\forall s}|}{|\mathcal{R}^{\forall r}|}$ and $\nu^{\mathsf{a}} : \text{RN} \to [0, 1)$.

Note that since $\exists r.\bot$ can never occur in any normalized $\mathcal{ALEH}$ concept description, we need not treat this case in Equation 7 (cf. Definition of homomorphism in the previous section and in (Baader and Küsters, 2006)). Intuitively, the homomorphism degree (hd) of the two given $\mathcal{ALEH}$ description trees can be computed based on the degree of common node label inclusion and the degree of common outgoing edges. Formula 6 calculates the proportion of the matched node labels comparing to all those available in the top level. Formula 7 and 9 computes the degrees of edge matching of an existential restriction and a universal restriction, respectively. If there is a shared edge label, then there is some degree of similarity; but the successors' labels and structures have yet to be checked. This is done recursively by calling the function hd$(\mathcal{G}_X, \mathcal{G}_Y)$.

The parameter $\mu^{\mathsf{e}}$ (resp., $\mu^{\mathsf{a}}$) defined in Formula 5 indicates how important the existentially quantified (resp., universally quantified) subconcepts are to be considered for similarity measure. The use of $\nu^{\mathsf{e}}$ and $\nu^{\mathsf{a}}$ allows to indicate an importance of the role name in an existential restriction and a universal restriction. It is similar to that in (Suntisrivaraporn, 2013) except that these are defined as a function on role names. This means that the importance of different role names and thus the discount of similarity between nested concepts can be unequally assigned based on their use and modelling discipline in a particular ontology. The value of $\gamma$ in Formula 8 and 10 indicates a degree of inclusion between the two edge labels. The case where $\gamma = 0$ means there is no commonality between two given roles, and hence further computation for the degrees of membership between their corresponding nested pairs should be omitted.

**Example.** To better understand how the algorithm works, consider the description tree $\mathcal{G}_{\text{Mother}}$ for the unfolding of Mother and the description tree

$\mathcal{G}_{\text{NonFosterFather}}$ for the unfolding of NonFosterFather as shown in Figure 2. Using $\mu$ as previously described and fixing $\nu^{\star}(r)$ to 0.4 for every role name $r \in \text{RN}$, the degrees of homomorphism from the root of $\mathcal{G}_{\text{Mother}}$ to the root of $\mathcal{G}_{\text{NonFosterFather}}$ can be computed as following steps (abbreviations are used for the sake of simplicity):

hd$(\mathcal{G}_{\text{M}}, \mathcal{G}_{\text{NFF}})$
$$:= \frac{2}{3}\text{p-hd}(\mathcal{P}_{\text{M}}, \mathcal{P}_{\text{NFF}}) + \frac{1}{3}\text{e-hd}(\mathcal{E}_{\text{M}}, \mathcal{E}_{\text{NFF}}) +$$
$$(0)\text{a-hd}(\mathcal{A}_{\text{M}}, \mathcal{A}_{\text{NFF}})$$
$$:= \frac{2}{3}[\frac{1}{2}] + \frac{1}{3}\text{e-hd}(\epsilon_i, \epsilon_j)$$
// with $\mu^{\mathsf{e}} = \frac{1}{3}$, $\mu^{\mathsf{a}} = 0$,
// $\epsilon_i = \exists\text{child.Person}$ and $\epsilon_j = \exists\text{child.Person}$
$$:= \frac{2}{3}[\frac{1}{2}] + \frac{1}{3}[\frac{1}{1}][\frac{2}{5} + \frac{3}{5}\text{hd}(\mathcal{G}_{\text{Person}}, \mathcal{G}_{\text{Person}})]$$
$$:= \frac{2}{3}[\frac{1}{2}] + \frac{1}{3}[\frac{2}{5} + \frac{3}{5}[1]]$$
$$:= \frac{2}{6} + \frac{1}{3}$$
$$:= 0.67$$

The reverse direction can be computed as follows:

hd$(\mathcal{G}_{\text{NFF}}, \mathcal{G}_{\text{M}})$
$$:= \frac{2}{4}\text{p-hd}(\mathcal{P}_{\text{NFF}}, \mathcal{P}_{\text{M}}) + \frac{1}{4}\text{e-hd}(\mathcal{E}_{\text{NFF}}, \mathcal{E}_{\text{M}}) +$$
$$\frac{1}{4}\text{a-hd}(\mathcal{A}_{\text{NFF}}, \mathcal{A}_{\text{M}})$$
$$:= \frac{2}{4}[\frac{1}{2}] + \frac{1}{4}\text{e-hd}(\epsilon_i, \epsilon_j) + \frac{1}{4}\text{a-hd}(\alpha_i, \alpha_j)$$
// with $\mu^{\mathsf{e}} = \frac{1}{4}$, $\mu^{\mathsf{a}} = \frac{1}{4}$,
// $\epsilon_i = \exists\text{child.Person}$ and $\epsilon_j = \exists\text{child.Person}$
// $\alpha_i = \forall\text{fchild.}\bot$ and $\alpha_j = \emptyset$
$$:= \frac{2}{4}[\frac{1}{2}] + \frac{1}{4}[\frac{1}{1}][\frac{2}{5} + \frac{3}{5}\text{hd}(\mathcal{G}_{\text{Person}}, \mathcal{G}_{\text{Person}})] + \frac{1}{4}[0]$$
$$:= \frac{1}{4} + \frac{1}{4}$$
$$:= 0.50$$

Hence, the degree of having a homomorphism from the root of $\mathcal{G}_{\text{Mother}}$ to $\mathcal{G}_{\text{NonFosterFather}}$ is 0.67, and that for the opposite direction is 0.50. The hd values for other pairs can be obtained in an analogous manner and are shown in Table 2.

Using a proof by induction, together with Theorem 1 (Baader and Küsters, 2000; Baader, 2003), it is

Table 2: Homomorphism degrees to and from the defined concepts in $\mathcal{O}_{\text{family}}$.

| hd($\downarrow$, $\rightarrow$) | Woman | Man | Parent | Mother | Father | MNS | MND | FF | NFF |
|---|---|---|---|---|---|---|---|---|---|
| Woman | 1.00 | 0.50 | 0.50 | 0.67 | 0.33 | 0.50 | 0.50 | 0.33 | 0.25 |
| Man | 0.50 | 1.00 | 0.50 | 0.33 | 0.67 | 0.25 | 0.25 | 0.67 | 0.50 |
| Parent | 0.50 | 0.50 | 1.00 | 0.67 | 0.67 | 0.43 | 0.43 | 0.50 | 0.50 |
| Mother | 1.00 | 0.50 | 1.00 | 1.00 | 0.67 | 0.68 | 0.68 | 0.50 | 0.50 |
| Father | 0.50 | 1.00 | 1.00 | 0.67 | 1.00 | 0.43 | 0.43 | 0.83 | 0.75 |
| MotherNoSon (MNS) | 1.00 | 0.50 | 1.00 | 1.00 | 0.67 | 1.00 | 0.85 | 0.50 | 0.60 |
| MotherNoDaughter (MND) | 1.00 | 0.50 | 1.00 | 1.00 | 0.67 | 0.85 | 1.00 | 0.50 | 0.60 |
| FosterFather (FF) | 0.50 | 1.00 | 1.00 | 0.67 | 1.00 | 0.43 | 0.43 | 1.00 | 0.75 |
| NonFosterFather (NFF) | 0.50 | 1.00 | 1.00 | 0.67 | 1.00 | 0.55 | 0.55 | 0.83 | 1.00 |

not difficult to obtain the correspondence between the homomorphism degree and subsumption.

**Proposition 2.** *Let $C, D$ be expanded and normalized $\mathcal{ALEH}$ concept descriptions, and $\mathcal{G}_C$, $\mathcal{G}_D$ be their corresponding description trees, respectively. Then, the following are equivalent:*

1. $C \sqsubseteq D$,
2. hd($\mathcal{G}_D$, $\mathcal{G}_C$) = 1.

In fact, the closer the hd($\mathcal{G}_D$, $\mathcal{G}_C$) value is to 1, the more likely the corresponding subsumption may hold. More precisely, the label and edge constraints in $\mathcal{G}_D$ can likely be simulated by those in $\mathcal{G}_C$.

# 4 $\mathcal{ALEH}$ SEMANTIC SIMILARITY

The homomorphism degree function introduced in Section 3 returns a degree that represents the similarity of one concept description compared to another concept description. As shown in the computation example, the direction of the homomorphism degree matters, viz., hd($\mathcal{G}_M$, $\mathcal{G}_{NFF}$) = 0.67, whereas hd($\mathcal{G}_{NFF}$, $\mathcal{G}_M$) = 0.50. Since both directions constitute the degree of the two concepts being equivalent, our similarity measure for $\mathcal{ALEH}$ concept descriptions is defined by means of these values.

**Definition 4** ($\mathcal{ALEH}$ Concept Similarity). *Let $C, D$ be expanded $\mathcal{ALEH}$ concept descriptions. The degree of similarity between $C$ and $D$ is defined as:*

$$\text{sim}_{\mathcal{ALEH}}(C, D) := \frac{\text{hd}(\mathcal{G}_C, \mathcal{G}_D) + \text{hd}(\mathcal{G}_D, \mathcal{G}_C)}{2} \quad (11)$$

Intuitively, the degree of similarity between two concepts is the average of the degree of having homomorphisms in both directions, thus $\text{sim}(C, D) = \text{sim}(D, C)$ as required.[1]

---

[1] Note that other functions apart from *average* could be applied; for instance, root mean square and multiplication (Suntisrivaraporn, 2013).

Based on the homomorphism degree values in Table 2, the degrees of similarity among the defined concepts in the example ontology $\mathcal{O}_{\text{family}}$ can be obtained; see Table 3. Note also that, though not included in Table 2 and 3, the similarity involving primitive concepts like Female and Person can also be computed. Nevertheless, the pairwise similarity degree between any two primitive concepts is zero by our definition since there is absolutely no commonality between them apart from both being subsumed by $\top$.

The similarity measure $\text{sim}_{\mathcal{ALEH}}$ generalizes sim for the DL $\mathcal{ELH}$ (Suntisrivaraporn, 2013; Tongphu and Suntisrivaraporn, 2015) in the sense that when two given concept descriptions are restricted to $\mathcal{ELH}$, then both measures coincide.

**Proposition 3.** *Let $C, D$ be two $\mathcal{ELH}$ concept descriptions. Then,*

$$\text{sim}_{\mathcal{ALEH}}(C, D) = \text{sim}(C, D).$$

This is the case since any $\mathcal{ELH}$ description tree is also an $\mathcal{ALEH}$ description tree that does not contain universal edges.

# 5 APPROXIMATING $\mathcal{ALCH}$ SEMANTIC SIMILARITY

A description logic $\mathcal{ALCH}$ can be considered as an extension of $\mathcal{ALEH}$ that supports more concept constructors, namely disjunction and full concept negation (see Table 1). Since DL $\mathcal{ALEH}$ is a language in the family DL $\mathcal{ALCH}$, in this section, we show that the notion of $\mathcal{ALEH}$ similarity measure can be extended to a new notion of $\mathcal{ALCH}$ similarity measure.

In Section 3 we review the structural characterization of subsumption $\mathcal{ALEH}$ through a homomorphism. Alas, this characterization is not directly applicable to the more expressive language $\mathcal{ALCH}$ due to disjunction. Fortunately, one can *approximate* an

Table 3: Similarity degree between a pair of defined concepts in $O_{\text{family}}$.

| hd($\downarrow$, $\rightarrow$) | Woman | Man | Parent | Mother | Father | MNS | MND | FF | NFF |
|---|---|---|---|---|---|---|---|---|---|
| Woman | 1.00 | 0.50 | 0.50 | 0.83 | 0.42 | 0.75 | 0.75 | 0.42 | 0.38 |
| Man | | 1.00 | 0.50 | 0.42 | 0.83 | 0.38 | 0.38 | 0.83 | 0.75 |
| Parent | | | 1.00 | 0.83 | 0.83 | 0.71 | 0.71 | 0.75 | 0.75 |
| Mother | | | | 1.00 | 0.67 | 0.84 | 0.84 | 0.58 | 0.58 |
| Father | | | | | 1.00 | 0.55 | 0.55 | 0.92 | 0.88 |
| MotherNoSon (MNS) | | | | | | 1.00 | 0.85 | 0.46 | 0.58 |
| MotherNoDaughter (MND) | | | | | | | 1.00 | 0.46 | 0.58 |
| FosterFather (FF) | | | | | | | | 1.00 | 0.79 |
| NonFosterFather (NFF) | | | | | | | | | 1.00 |

$\mathcal{ALCH}$-concept description in the less expressive DL $\mathcal{ALEH}$. Once approximation is calculated, the similarity measure introduced in this paper could be used to obtain approximate similarity between two concept descriptions written in $\mathcal{ALCH}$.

**Definition 5** (Approximation). *(Baader and Küsters, 2006) Let C be an $\mathcal{ALCH}$-concept description. An $\mathcal{ALEH}$-concept description D is an $\mathcal{ALEH}$-approximation of C, written $\mathcal{ALEH}$-approx(C), iff*

  i. $C \sqsubseteq D$ *and*
  ii. $D \sqsubseteq E$ *for every $\mathcal{ALEH}$-concept description E with $C \sqsubseteq E$.*

Intuitively, an approximation is the most specific concept in $\mathcal{ALEH}$ that subsumes the given $\mathcal{ALCH}$ concept. One can approximate an $\mathcal{ALCH}$ concept by resorting to finding commonalities among sub-concepts in a disjunction, also known as the *least common subsumer (lcs)* problem (Turhan, 2007; Baader et al., 1998).

We define the notion of similarity measure between two $\mathcal{ALCH}$ concept descriptions as follows:

**Definition 6** ($\mathcal{ALCH}$ Concept Similarity). *Let X,Y be $\mathcal{ALCH}$ concept descriptions. The degree of similarity between X and Y, in symbols $\text{sim}_{\mathcal{ALCH}}(X,Y)$, is defined as:*

$$\text{sim}_{\mathcal{ALCH}}(X,Y) := \\ \text{sim}_{\mathcal{ALEH}}(\mathcal{ALEH}\text{-approx}(X), \mathcal{ALEH}\text{-approx}(Y))$$

An analogous idea can be employed to compute concept similarity in another DLs and yet using another similarity measure. For instance, it is possible to approximate $\mathcal{ELU}$-concept descriptions ($\mathcal{EL}$ extended with disjunction) and then compute similarity using the known measure for $\mathcal{EL}$ (Lehmann and Turhan, 2012; Suntisrivaraporn, 2013). It remains however to be shown whether this produces acceptable similarity results in practice.

# 6 RELATED WORKS

The subject of concept similarity has been widely studied. The techniques can be roughly classified into two main groups: a structure-based approach and an edit-distance-based approach.

In (Distel et al., 2014), the authors introduced a new framework of concept similarity measure. This framework is based on a counting of relaxation operations. A similarity is defined by means of the distance between concept descriptions C and D, i.e. the number of times D needs to be relaxed before it subsumes C. The method is claimed to satisfied several properties of concept similarity but has not yet been implemented.

A measure proposed by (Ge and Qiu, 2008) calculates a degree of similarity based on the depth of a concept defined in different levels of the ontological hierarchy. The method considers the distance relationship (subsumption relation) between concepts and assigned different weights to the role depth. The degree of similarity between two concepts was measured by means of a distance (a propagation of all label weights) to their least common subsumer. Similar approaches were proposed in (Ge and Qiu, 2003; Giunchiglia et al., 2007). Despite their usefulness in structural analysis, these methods were fully relied on an ontology hierarchy and usually ignored constraints of concept definitions in the ontology.

A simple method for similarity measure in the basic DL $\mathcal{L}_0$ (i.e. no use of roles) was proposed in (Jaccard, 1901), known as *Jaccard Index*. An extension thereof to the DL $\mathcal{ELH}$ was proposed in (Lehmann and Turhan, 2012). The extended work suggested a new framework that satisfies several properties for similarity. While the framework is defined in general, the functions and operators needed for the computation are parameterized and thus left to be specified. Moreover, the framework does not contain implementation details.

The notion of homomorphism degree was originally introduced in (Suntisrivaraporn, 2013) and employed as the heart of the similarity measure for the DL $\mathcal{EL}$. This has been extended to $\mathcal{ELH}$ and continuously studied in (Tongphu and Suntisrivaraporn, 2014; Tongphu and Suntisrivaraporn, 2015).

Racharak and Suntisrivaraporn suggested two new notions of similarity for the DL $\mathcal{FL}_0$ (Racharak and Suntisrivaraporn, 2015). Both the skeptical and credulous similarity measures are derived from the known structural characterization subsumption through inclusion of regular languages.

The similarity measure presented in this paper is similar to those reported in (Tongphu and Suntisrivaraporn, 2014; Suntisrivaraporn, 2013). It however focuses on the strictly more expressive DL and employs generalizations of the normalization and characterization from (Baader and Küsters, 2006).

## 7 DISCUSSIONS AND FUTURE WORKS

This paper presents a new notion of concept similarity for the DL $\mathcal{ALEH}$ w.r.t. an unfoldable terminology and suggests a way to approximate concept similarity for the more expressive $\mathcal{ALCH}$. At the heart of the measure is the calculation of the degree of homomorphism to and from between two description trees. To allow this, we first review and extend the known normalization and homomorphism to take into account also role hierarchy axioms. The proposed similarity measure can be regarded as an extension of the similarity measure sim for the $\mathcal{EL}$ family (Suntisrivaraporn, 2013; Tongphu and Suntisrivaraporn, 2015).

There are various directions for future works. One could try to evaluate the proposed measure on appropriate ontologies from real-world domains. Similar to the experiments on SNOMED CT reported in (Tongphu and Suntisrivaraporn, 2015), a similar setting can be carried out. Besides, more expressive ontologies that make use of the universal quantification such as GALEN could be experimented upon. It can be expected to find out new hidden knowledge in the ontology that could not have been done before with the mere standard reasoner. Another useful application is a measure of similarity between diseases proposed in (Mathur and Dinakarpandian, 2012). The application has shown useful cases in similarity measure processes underlying each disease for more accurate unknown disease prediction.

Concerning the choice of representation language, it is an obvious future work to explore non-approximate similarity measure for $\mathcal{ALC}$ by investigating under scrutiny into the original tableau algorithm. Another direction for future work could be to compare the measure presented in this paper to those two notions of similarity for $\mathcal{FL}_0$ introduced in (Racharak and Suntisrivaraporn, 2015). Since $\mathcal{FL}_0$ is a sub-logic of $\mathcal{ALEH}$ and as such $\text{sim}_{\mathcal{ALEH}}$ is applicable also to $\mathcal{FL}_0$, it is interesting to explore whether $\text{sim}_{\mathcal{ALEH}}$ is *stronger* (see (Racharak and Suntisrivaraporn, 2015)) than the skeptical and credulous similarity measures.

## ACKNOWLEDGEMENTS

## REFERENCES

Baader, F. (2003). Terminological cycles in a description logic with existential restrictions. In Gottlob, G. and Walsh, T., editors, *Proceedings of the 18th International Joint Conference on Artificial Intelligence*, pages 325–330. Morgan Kaufmann.

Baader, F., Calvanese, D., McGuinness, D., Nardi, D., and Patel-Schneider, P., editors (2007). *The Description Logic Handbook: Theory, Implementation and Applications*. Cambridge University Press, second edition.

Baader, F. and Küsters, R. (2000). Matching in description logics with existential restrictions. In A.G. Cohn, F. Giunchiglia, and B. Selman, editors, *Proceedings of the Seventh International Conference on Knowledge Representation and Reasoning (KR2000)*, pages 261–272, San Francisco, CA. Morgan Kaufmann Publishers.

Baader, F. and Küsters, R. (2006). Nonstandard inferences in description logics: The story so far. In Gabbay, D., Goncharov, S., and Zakharyaschev, M., editors, *Mathematical Problems from Applied Logic I*, volume 4 of *International Mathematical Series*, pages 1–75. Springer-Verlag.

Baader, F., Küsters, R., and Molitor, R. (1998). Computing least common subsumers in Description Logics with existential restrictions. LTCS-Report LTCS-98-09, LuFG Theoretical Computer Science, RWTH Aachen, Germany. See http://www-lti.informatik.rwth-aachen.de/Forschung/Papers.html.

Distel, F., Atif, J., and Bloch, I. (2014). Concept dissimilarity with triangle inequality. In *Proceedings of the Fourteenth International Conference on Principles of*

*Knowledge Representation and Reasoning (KR'14)*, Vienna, Austria. AAAI Press. Short Paper. To appear.

Ge, J. and Qiu, Y. (2003). Concept similarity matching based on semantic distance. In Gottlob, G. and Walsh, T., editors, *Proceedings of the Forth International Conference on Semantics, Knowledge and Grid (SKG 2008)*, pages 380–383. Morgan Kaufmann.

Ge, J. and Qiu, Y. (2008). Concept similarity matching based on semantic distance. In *SKG*, pages 380–383. IEEE Computer Society.

Giunchiglia, F., Yatskevich, M., and Shvaiko, P. (2007). Semantic matching: Algorithms and implementation. *Journal of Data Semantics*, 9:1–38.

Jaccard, P. (1901). Étude comparative de la distribution florale dans une portion des Alpes et des Jura. *Bulletin del la Société Vaudoise des Sciences Naturelles*, 37:547–579.

Lehmann, K. and Turhan, A.-Y. (2012). A framework for semantic-based similarity measures for $\mathcal{ELH}$ - concepts. In del Cerro, L. F., Herzig, A., and Mengin, J., editors, *JELIA*, volume 7519 of *Lecture Notes in Computer Science*, pages 307–319. Springer.

Mathur, S. and Dinakarpandian, D. (2012). Finding disease similarity based on implicit semantic similarity. *Journal of Biomedical Informatics*, 45(2):363–371.

Racharak, T. and Suntisrivaraporn, B. (2015). Similarity measures for $\mathcal{FL}_0$ concept descriptions from an automata-theoretic point of view. In *Information and Communication Technology for Embedded Systems (IC-ICTES)*, pages 1–6. IEEE Computer Society.

Suntisrivaraporn, B. (2013). A similarity measure for the description logic $\mathcal{EL}$ with unfoldable terminologies. In *International Conference on Intelligent Networking and Collaborative Systems (INCoS-13)*, pages 408–413.

Tongphu, S. and Suntisrivaraporn, B. (2014). On desirable properties of the structural subsumption-based similarity measure. In *Semantic Technology - 4th Joint International Conference, JIST 2014, Chiang Mai, Thailand, November 9-11, 2014. Revised Selected Papers*, pages 19–32.

Tongphu, S. and Suntisrivaraporn, B. (2015). Algorithms for measuring similarity between elh concept descriptions: A case study on SNOMED CT. *Journal of Computing and Informatics*. (Accepted in May 2015; To appear).

Turhan, A.-Y. (2007). *On the Computation of Common Subsumers in Description Logics*. PhD thesis, TU Dresden, Institute for Theoretical Computer Science, Germany.