# Modelling OLIF frame with EAGLES/ISLE specifications:
# an interlingual approach

## *El modelado de OLIF utilizando las especificaciones de EAGLES/ISLE:*
## *un enfoque interlingüístico*

**Carlos Periñán-Pascual,  Francisco Arcas-Túnez**
Universidad Católica San Antonio
Campus de los Jerónimos s/n
30107 Guadalupe - Murcia (Spain)
{jcperinan, farcas}@pdi.ucam.edu

**Resumen**: FunGramKB es una base de conocimiento léxico-conceptual para su implementación en sistemas del PLN. El modelo léxico de FunGramKB se construyó a partir del modelo de OLIF, aunque fue preciso incorporar algunas de las recomendaciones de EAGLES/ISLE con el fin de poder diseñar lexicones computacionales más robustos. El propósito de este artículo es describir cómo el enfoque interlingüístico de FunGramKB proporciona una visión más cognitiva de los marcos léxicos que las propuestas por OLIF y EAGLES/ISLE.
**Palabras clave**: FunGramKB, OLIF, EAGLES, ISLE, lexicón, ontología, marco, postulado de significado.

**Abstract**: FunGramKB is a lexico-conceptual knowledge base for NLP systems. The FunGramKB lexical model is basically derived from OLIF and enhanced with EAGLES/ISLE recommendations with the purpose of designing robust computational lexica. However, the FunGramKB interlingual approach gives a more cognitive view to EAGLES/ISLE proposals. The aim of this paper is to describe how this approach influences the way of conceiving lexical frames.
**Keywords**: FunGramKB, OLIF, EAGLES, ISLE, lexicon, ontology, frame, meaning postulate.

## 1 Introduction

FunGramKB (Functional Grammar Knowledge Base) is a lexico-conceptual knowledge base for NLP systems, mainly those requiring natural language understanding. FunGramKB is multipurpose, in the sense that it is both multifunctional and multilanguage. In other words, FunGramKB has been designed to be reused in various NLP tasks (e.g. information retrieval/extraction, machine translation or dialogue-based systems) and with several natural languages.[1]

The FunGramKB lexical model is basically derived from OLIF[2] (Lieske et al. 2001; McCormick 2002; McCormick et al. 2004) and

---

[1] FunGramKB lexica for English and Spanish are being currently populated.

[2] OLIF (Open Lexicon Interchange Format) is created in the 90's as part of the OTELO (Open Translation Environment for Localization) project, whose primary goal is the development of interfaces and formats which can help users share lexical resources within the translation environment (e.g. machine translation, translation memories, terminology databases, and so on).

enhanced with EAGLES/ISLE[3] recommendations (EAGLES 1993, 1996a, 1996b, 1999; Monachini et alii 2003; Underwood and Navarretta 1997; Calzolari et alii 2001a, 2001b, 2003). OLIF, an XML-compliant standard for lexical/terminological data encoding, was chosen as the starting point for implementing the FunGramKB lexical level. However, some parts of the OLIF model had to be re-considered in order to make it conform to the FunGramKB architecture.[4] The FunGramKB team soon realised that, for example, full-fledged lexical frames were not possible if language engineers were confined to OLIF recommendations. Therefore, OLIF was modelled with EAGLES/ISLE specifications with the purpose of designing robust computational lexica.

In computational linguistics, lexical frames usually include key information which allows the computer to build the underlying predication of an input text. This paper presents a conceptualist model of frame semantics which, in turn, complies with current standards for computational lexica. Section 2 briefly describes the two-tier architecture of the FunGramKB model. Section 3 shows how frame participants should be fully integrated into the lexical meaning of verbs via meaning postulates, resulting in a more "intelligent" resource for natural language understanding. Finally, sections 4 and 5 discuss the degree to which FunGramKB is indebted to OLIF and EAGLES/ISLE standards.

---

[3] EAGLES (The Expert Advisory Group on Language Engineering Standards) is an initiative sponsored by the European Commission which aims to provide recommendations for the standardization of the language technologies field. More particularly, the Computational Lexicons Interest Group is in charge of analysing the main practices in lexicographic encoding by comparing computational lexical resources available in European languages.

ISLE (International Standards for Language Engineering) is initiated in 2000 as an extension of EAGLES work. The objective of this joint EU-US project is to support R&D on Human Language Technology issues. The ISLE Computational Lexicon Working Group is committed to the design of MILE (Multilingual ISLE Lexical Entry), a meta-entry for the encoding of multilingual lexical information.

[4] Indeed, one of the advantages of OLIF is the ease of extensibility and customization of its XML-based format in order to accommodate it to the requirements of a project.

## 2 The FunGramKB architecture

FunGramKB comprises two information levels, where several independent modules are interrelated:[5]

Lexical level (i.e. linguistic knowledge):
- The lexicon stores morphosyntactic, pragmatic and collocational information of lexical units.
- The morphicon helps our system to handle cases of inflectional morphology.

Cognitive level (i.e. non-linguistic knowledge):
- The ontology is presented as a hierarchical structure of well-defined concepts used by ordinary humans when talking about everyday situations.
- The cognicon stores procedural knowledge by means of cognitive macrostructures, i.e. script-like schemata in which a sequence of stereotypical actions is organised on the basis of temporal continuity.
- The onomasticon stores information about instances of entities, such as people, cities, products, and so on.

The motivation of this two-tier design lies in the fact that lexical modules are language specific but cognitive modules are shared by all languages. In other words, computational lexicographers must develop one lexicon and one morphicon for English, one lexicon and one morphicon for Spanish and so on, but knowledge engineers build just one ontology, one cognicon and one onomasticon to process any language input cognitively. Unlike most current NLP systems, where the lexicalist approach prevails, the FunGramKB architecture is ontology-oriented, since the ontology plays a pivotal role between the lexical and the cognitive levels.

---

[5] Computationally speaking, entries for any of these modules take the form of XML-formatted data structures. XML was chosen as the formal language for knowledge representation because data can be encoded in such a portable way that information can be easily compilable into the format that is needed by other formalisms and systems.
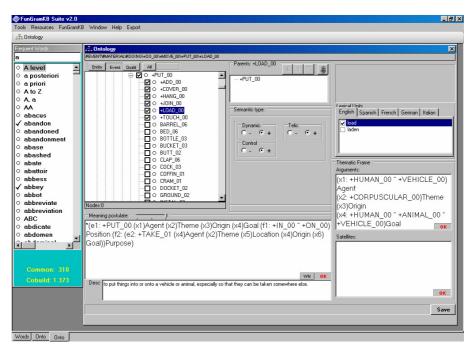
Figure 1: FunGramKB Suite

Since FunGramKB is intended to be an extensive NLP knowledge base, it is important to minimize the effort for ontology maintenance, so strict control is placed on the management of data consistency. As shown in figure 1, FunGramKB Suite has been designed for that purpose.

For instance, the construction of knowledge schemata such as predicate frames or meaning postulates is semiautomatic, because human intervention is required but the knowledge engineer's intuition is guided and reviewed through FunGramKB Editor, so that consistent well-formed constructs can be stored.

The following section describes how the FunGramKB conceptualist approach undoubtedly influences the way of conceiving frames.
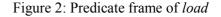
## 3 Integrating frames into meaning postulates

Most semantic representations of verbs have traditionally taken one of two forms (Levin 1995): semantic role-centred approaches (Fillmore 1968, Gruber 1965), where verb arguments are identified on the basis of their semantic relations with the verb, or predicate decomposition approaches (Jackendoff 1972, Schank 1973), which involve the decomposition of verb meaning by means of a restricted set of primitive predicates.

In FunGramKB, both approaches are integrated. Similarly to semantic role-centred approaches, verbs are assigned one or more

frames, which are called "predicate frames". To illustrate, figure 2 displays both the parenthetic string representation (edition format) and the XML representation (storage format) of the predicate frame of *load*:

$(x1)_{S/Agent/NP}$ $(x2)_{O/Theme/NP}$ $(x4)_{A/Goal/PP-into,onto}$
$(x1)_{S/Agent/NP}$ $(x4)_{O/Goal/NP}$ $(x2)_{A/Theme/PP-with}$

```
<prFrame>
    <pattern>
        <slot phrase="NP" syn="S" sem="Agent" />
            <slot phrase="NP" syn="O" sem="Theme" />
            <slot phrase="PP" syn="A" sem="Goal">
                <prep>into</prep>
                <prep>onto</prep>
            </slot>
    </pattern>
    <pattern>
        <slot phrase="NP" syn="S" sem="Agent" />
        <slot phrase="NP" syn="O" sem="Goal" />
            <slot phrase="PP" syn="A" sem="Theme">
                <prep>with</prep>
            </slot>
    </pattern>
</prFrame>
```

Figure 2: Predicate frame of *load*

The predicate frame is a structural scheme in which the quantitative and qualitative[6]

---

[6] Selectional preferences on an argument are not really stored in predicate frames, but they are part of thematic frames in the FunGramKB ontology. However, since predicate frames are derived from thematic frames, selectional preferences can definitely take part in full-fledged predicate frames.

valencies of the verb are stated: e.g. *load* has three subcategorized arguments with the semantic functions Agent, Theme and Goal. Moreover, predicate frames are enriched with information about subcategorization patterns describing the phrasal realizations and syntactic behaviour of the arguments which can linguistically co-occur with the verb.

On the other hand, and like predicate decomposition approaches, a lexical unit is linked to a meaning postulate through a conceptual unit in the FunGramKB ontology.[7] Furthermore, predicate frames assigned to a lexical unit are integrated into the meaning representation to which the lexical unit is linked by means of the "thematic frame". To illustrate, figure 3 displays both the parenthetic string representation and the XML representation of the thematic frame of +LOAD_00:

$(x1: +HUMAN\_00 \,^\wedge\, +VEHICLE\_00)_{Agent}$ $(x2: +CORPUSCULAR\_00)_{Theme}$
$(x3)_{Origin}$ $(x4: +HUMAN\_00 \,^\wedge\, +ANIMAL\_00 \,^\wedge\, +VEHICLE\_00)_{Goal}$

```
<thFrame>
    <Arguments>
        <x n="1" sem="Agent">
            <PrefSet oper="xor">
                <Pref concept="+HUMAN_00" />
                <Pref concept="+VEHICLE_00" />
            </PrefSet>
        </x>
        <x n="2" sem="Theme">
            <Pref concept="+CORPUSCULAR_00" />
        </x>
        <x n="3" sem="Origin" />
        <x n="4" sem="Goal">
            <PrefSet oper="xor">
                <Pref concept="+HUMAN_00" />
                <Pref concept="+ANIMAL_00" />
                <Pref concept="+VEHICLE_00" />
            </PrefSet>
        </x>
    </Arguments>
</thFrame>
```

Figure 3: Thematic frame of +LOAD_00

Thematic frames are cognitive schemata specifying the type of participants involved in the situation described by an event. These participants can be instantiated in the form of arguments in the predicate frames assigned to the lexical units linked to that event.[8] Therefore, predicate frames are lexical constructs belonging to a particular language, but they are constructed from the interlingual thematic frames located in the ontology. In FunGramKB, every argument found in the predicate frame of a verb must be referenced through co-indexation in the thematic frame of the event to which the verb is linked. Moreover, every argument found in the thematic frame of an event is referenced through co-indexation in the meaning postulate assigned to that event. To illustrate, figure 4 displays both the parenthetic string representation and the XML representation of the meaning postulate of +LOAD_00:

$+(e1: +PUT\_00 \; (x1)_{Agent} \; (x2)_{Theme} \; (x3)_{Origin} \; (x4)_{Goal} \; (f1: +IN\_00 \,^\wedge\, +ON\_00)_{Position} \; (f2: (e2: +TAKE\_01 \; (x4)_{Agent} \; (x2)_{Theme} \; (x5)_{Location} \; (x4)_{Origin} \; (x6)_{Goal}))_{Purpose})$

```
<mPostulate>
    <Predication opr="+">
        <e n="1" concept="+PUT_00">
            <Arguments>
                <x n="1" sem="Agent" />
                <x n="2" sem="Theme" />
                <x n="3" sem="Origin" />
                <x n="4" sem="Goal" />
            </Arguments>
            <Satellites>
                <fSet oper="and">
                    <f n="1" sem="Position">
                        <PrefSet oper="xor">
                            <Pref concept="+IN_00" />
                            <Pref concept="+ON_00" />
                        </PrefSet>
                    </f>
                    <f n="2" sem="Purpose">
                        <e n="2" concept="+TAKE_01">
                            <Arguments>
                                <x n="4" sem="Agent" />
                                <x n="2" sem="Theme" />
                                <x n="5" sem="Location" />
                                <x n="4" sem="Origin" />
                                <x n="6" sem="Goal" />
                            </Arguments>
                        </e>
                    </f>
                </fSet>
            </Satellites>
        </e>
    </Predication>
</mPostulate>
```

Figure 4: Meaning postulate of +LOAD_00

---

[7] In fact, regularities in the semantic distribution of verbs in FunGramKB are not based on syntactic criteria (cf. Levin 1993) but on the cognitive decompositions of events by means of their meaning postulates.

[8] The difference between thematic frames and predicate frames is partly influenced by the distinction in the Construction Grammar (Goldberg 1995) between argument roles and participant roles respectively, where the first are related to the construction and the latter to the frame of a particular verb.

For example, the first predicate frame of *load* matches the morphosyntactic structure of a sentence such as *They loaded all their equipment into backpacks*, identifying *they* as the loaders (Agent), *equipment* as the thing to be loaded (Theme) and *backpacks* as the target entity where that thing is placed (Goal). However, the semantic burden of the frame is greater when linked to the thematic frame and the meaning postulate of +LOAD_00, which reveal that "they put the equipment into backpacks because they intended to carry it to another place".[9]

As it has been demonstrated, every argument in the predicate frame of a verb is finally integrated in the meaning postulate of its event through the arguments of its thematic frame, which plays a crucial role in both the semantic role-centred and predicate decomposition approaches to the semantic representation of verbs in FunGramKB.

## 4 The OLIF frame category

Three OLIF data categories are relevant for the construction of FunGramKB predicate frames:

(i) &lt;transType&gt; specifies the type of prototypical transitivity of the verb.

(ii) &lt;synFrame&gt; describes the subcategorization of the lexical entry. A slot-grammar approach is taken for the description of syntactic frames. For example, the frame for the English verb *try* is as follows (McCormick 2002):

[subj, (dobj-opt | dobj-sent-ing-opt | dobj-sent-inf-opt)]

(iii) &lt;prep&gt; specifies the preposition that fills a "prepositional phrase" slot.

The main advantage of the FunGramKB model of predicate frame does not lie just on the further specification of the lexical

---

[9] Indeed, a lexical unit is associated to much more semantic information which is really shown in its meaning postulate. In FunGramKB, all this underlying cognitive information is revealed through a multi-level process called MicroKnowing (Periñán-Pascual and Arcas-Túnez 2005), where thematic frames also play a key role in the application of the inheritance and inference mechanisms on meaning postulates.

information, but also on its remarkable conceptualist approach. To this respect, two main differences are observed between OLIF frames and FunGramKB predicate frames. Firstly, OLIF frames are semantically underspecified, since no semantic role is assigned to any slot. Secondly, slot fillers in OLIF are language-specific and not formally represented, whereas in FunGramKB selectional preferences are represented by concepts. Selection preferences should not be lexicalized, but somehow they should be part of human beings' cognitive knowledge. The benefit of this approach is twofold: (i) the use of concepts as the building blocks of predicate frames removes the problem of lexical semantic ambiguity, and (ii) the inferential power of the reasoning engine is more robust if predictions are based on cognitive expectations. The following section highlights the influence of EAGLES/ISLE standard on the construction of both predicate and thematic frames in FunGramKB.

## 5 Taking into account EAGLES/ISLE recommendations

EAGLES/ISLE proposes two types of frame: the syntactic frame, which describes the surface structure, and the semantic frame, which describes the deep structure.

On the one hand, the syntactic (or subcategorization) frame is expressed as a list of slots, where each slot is described in terms of phrasal realization, grammatical function, restricting features and optionality. Indeed, EAGLES/ISLE proposes a FrameSet to be included in the syntactic entry with the aim of collecting surface regular alternations associated with the same deep structure by explicitly linking the slots of the alternating frames by means of rules. Frames involved in a FrameSet are considered to be at the same level, i.e. no alternating frame has a status of privilege from which the other frames are derived through some lexical rule. Surprisingly, the EAGLES/ISLE approach is not as descriptively economical as the traditional approach, where, given two alternating frames, one of them is deemed to be basic and the other derivative.

In comparison with the EAGLES/ISLE proposal of syntactic frame, FunGramKB predicate frames make a limited use of restricting features, because only lexical features can be used to refine the information

specified in the arguments: e.g. the preposition that introduces a prepositional phrase. Moreover, the optional realization of an argument is not stated in FunGramKB predicate frames, because it is thought that context can admit the omission of any traditionally obligatory argument. Concerning frame alternations, FunGramKB can reflect all those syntactic phenomena in which no satellite is involved in the shift. On the contrary, satellite-oriented alternations such as locative alternations or material/product alternations are disregarded, since satellites are excluded from predicate frames.

On the other hand, the EAGLES/ISLE semantic frame (or argument structure) is defined in the form of a predicate and a list of arguments, which are described in terms of thematic role and semantic preferences. In general, the type of information in the FunGramKB thematic frame matches that of the EAGLES/ISLE semantic frame; however, differences are found in their approaches to the syntax-semantics interface within a multilingual dimension. EAGLES/ISLE recommends preferably a transfer architecture,[10] where monolingual syntactic and semantic frames are put into correlation between L1 and L2; in addition, this approach requires the specification of a set of transformational operations to go from L1 to L2. On the contrary, an interlingual model is adopted by FunGramKB, where thematic frames serve as the bridge between L1 predicate frames and those in L2. Transfer rules are not required since thematic frames are not linked to any particular lexicon but to the ontology, which is shared by all languages.

As a result, the FunGramKB interlingual approach gives a more cognitive view to the EAGLES/ISLE semantic frame. Firstly, EAGLES/ISLE recommends that both the predicate and its arguments should be instantiated with language-dependent lexical units, so that complexity in the linkage of the syntactic and semantic frames is dramatically reduced. On the contrary, sub-elements in FunGramKB thematic frames are not lexically driven, since predicates and semantic preferences on arguments are chosen from concepts of the ontology. Therefore, the notion

of thematic frame is more abstract than that of semantic frame. Secondly, EAGLES/ISLE proposes that the choice of the number of arguments for a predicate should be determined on purely semantic grounds; thus it is possible that (a) a syntactic position cannot be mapped to any semantic argument—i.e. reduced correspondence, or (b) a semantic argument cannot be mapped to any syntactic position—i.e. augmented correspondence. In FunGramKB, any decision on the type and number of arguments in thematic frames is guided by cognitive criteria. However, the FunGramKB architecture is so marked by the conceptualist approach that, for example, reduced correspondences in the syntax-semantics interface are not permitted because predicate frames are built out of their thematic frames, but not conversely.

## 6 Conclusions and future work

This paper presents the modifications and extensions to the OLIF model of frame by taking into account some of the EAGLES/ISLE recommendations. The result is that FunGramKB is provided with predicate frames in the lexicon (lexical frames) and thematic frames in the ontology (cognitive frames). We have also described that the two most important approaches to lexical semantic representation are fully integrated in FunGramKB: thus verbs are assigned one or more predicate frames, whose arguments play an active role in the construction of the meaning postulates to which those verbs are linked. In short, the FunGramKB interlingual approach, which gives a more cognitive view to the EAGLES/ISLE semantic frame, contributes to the large-scale development of deep-semantic NLP resources, mainly for natural language understanding.

We intend to develop a more robust characterization of predicate frames by exploring linguistically annotated corpora. Thus, and guided by some other suggestions proposed by EAGLES/ISLE, predicate frames could also include:

(i) an index indicating the frequency of the frame,[11]

---

[10] Although other approaches to translation are also considered, EAGLES/ISLE multilingual layer is inspired mostly on the transfer-based model.

[11] Frame probability can be particularly useful in natural language generation. For example, the current model of FunGramKB stores a default translation equivalent for every lexical unit, but it could be possible to use statistical information to

(ii)    a wider range of participants, i.e. satellites together with arguments,

(iii)    morphosyntactic restrictions on participants, e.g. whether the phrasal realization in a slot must be instantiated via plural word form,

(iv)    conditional optionality of participants, i.e. when the absence of a participant excludes or requires the presence of another participant,

(v)    lexical collocations as selectional preferences on participants,

### Bibliography

Calzolari, N., R. Grishman, and M. Palmer. eds. 2001. Survey of major approaches towards bilingual/multilingual lexicons. ISLE Deliverable D2.1-D3.1. ISLE Computational Lexicon Working Group.

Calzolari, N., F. Bertagna, A. Lenci, and M. Monachini. eds. 2003. Standards and best practice for multilingual computational lexicons and MILE. Deliverable D2.2-D3.2. ISLE Computational Lexicon Working Group.

Calzolari, N., A. Lenci, and A. Zampolli. 2001a. The EAGLES/ISLE computational lexicon working group for multilingual computational lexicons. *Proceedings of the First International Workshop on Multimedia Annotation*. Tokyo (Japan).

Calzolari, N., A. Lenci, and A. Zampolli. 2001b. International standards for multilingual resource sharing: the ISLE Computational Lexicon Working Group. *Proceedings of the ACL 2001 Workshop on Human Language Technology and Knowledge Management*. 71-78, Morristown (USA).

EAGLES Lexicon Interest Group. 1993. EAGLES: Computational Lexicons Methodology Task. EAGLES Document EAG-CLWG-METHOD/B.

EAGLES Lexicon Interest Group. 1996a. EAGLES: synopsis and comparison of morphosyntactic phenomena encoded in lexicons and corpora. A common proposal and applications to European languages. EAGLES Document EAG-CLWG-MORPHSYN/R.

EAGLES Lexicon Interest Group. 1996b. EAGLES: preliminary recommendations on subcategorisation. EAGLES Document EAG-CLWG-SYNLEX/P.

EAGLES Lexicon Interest Group. 1999. EAGLES: preliminary recommendations on lexical semantic encoding. Final report LE3-4244.

Fillmore, C.J. 1968. The case for case. E. Bach and R.T. Harms. eds. *Universals in Linguistic Theory*. Holt, Rinehart & Winston, New York, 1-88.

Goldberg, A.E. 1995. *Constructions: A Construction Grammar Approach to Argument Structure*. The University of Chicago Press, Chicago.

Gruber, J.S. 1965. *Studies in Lexical Relations*. Doctoral dissertation. MIT.

Jackendoff, R.S. 1972. *Semantic Interpretation in Generative Grammar*. MIT Press, Cambridge (Mass.).

Levin, B. 1993. *English Verb Classes and Alternations: A Preliminary Investigation*. The University of Chicago Press, Chicago.

Levin, B. 1995. Approaches to lexical semantic representation. D.E. Walker, A. Zampolli, and N. Calzolari. eds. *Automating the Lexicon: Research and Practice in a Multilingual Environment*. Oxford University Press, New York.

Lieske, C., S. McCormick, and G. Thurmair. 2001. The Open Lexicon Interchange Format (OLIF) comes of age. *Proceedings of the Machine Translation Summit VIII: Machine Translation in the Information Age*. 211-216, Santiago de Compostela (Spain).

McCormick, S. 2002. *The Structure and Content of the Body of an OLIF v.2.0/2.1*. The OLIF2 Consortium.

McCormick, S., C. Lieske, and A. Culum. 2004. *OLIF v.2: A Flexible Language Data Standard*. The OLIF2 Consortium.

Monachini, M., F. Bertagna, N. Calzolari, N. Underwood, and C. Navarretta. 2003. *Towards a Standard for the Creation of*

address the translation of an L1 lexical unit to the most probable equivalent in L2.

*Lexica*. ELRA European Language Resources Association.

Periñán-Pascual, C. and F. Arcas-Túnez. 2005. Microconceptual-Knowledge Spreading in FunGramKB. *9th IASTED International Conference on Artificial Intelligence and Soft Computing*, 239- 244, ACTA Press, Anaheim-Calgary-Zurich.

Schank, R.C. 1973. Identification of conceptualizations underlying natural language. R.C. Schank and K.M. Colby. eds. *Computer Models of Thought and Language*. W.H. Freeman, San Francisco, 187-247.

Underwood, N. and C. Navarretta. 1997. Towards a standard for the creation of lexica. Center for Sprogteknologi. Copenhagen.