

# OAL: una herramienta para el desarrollo de recursos lingüísticos

## *OAL: a tool for the development of linguistic resources*

Javier Couto<sup>1,2</sup>, Helena Blancafort<sup>1,3</sup>, Somara Seng<sup>1</sup>, Anass Talby<sup>1</sup>, Claude de Loupy<sup>1,2</sup>

<sup>1</sup> Syllabs  
Paris, France  
{couto,blancafort,loupy}@syllabs.com

<sup>2</sup> Laboratoire Modyco  
Université de Paris 10  
Nanterre, France

<sup>3</sup> Universitat Pompeu Fabra  
Barcelona, España

**Resumen:** OAL es un conjunto de herramientas para ayudar a los lingüistas a crear, alimentar y gestionar léxicos computacionales. Está diseñado como un cliente-servidor que incluye un agente inteligente de aspiración de contenido textual Web y una cadena para trabajar a partir de corpus bruto.

**Palabras clave:** editor de léxicos, aspirador de contenido textual Web.

**Abstract:** OAL is a tool to assist the linguist in the development, enrichment and management of computational lexicons. It has been designed as a client-server which includes an intelligent web crawler and a chain to process raw corpora.

**Keywords:** lexicon editor, web crawler.

## 1 Introducción

Uno de los aspectos claves en el procesamiento de lenguaje natural está vinculado con la calidad de los recursos lingüísticos (p.ej. léxicos morfosintácticos) sobre los cuales se basan las aplicaciones. Aún hoy, su creación y enriquecimiento suponen un proceso costoso y, muchas veces, tedioso.

Sin herramientas específicas, quien desarrolla estos recursos necesita mancomunarse los resultados de distintos tratamientos: aspiradores de contenido Web, *scripts* para lidiar con problemas de codificado e idioma, *scripts* para extraer el contenido textual, para realizar importaciones y exportaciones de información, entre otros. La información suele guardarse en formato textual y la manipulación de varios ficheros es una práctica corriente.

Para resolver esta problemática, Syllabs dispone de OAL<sup>1</sup>, un conjunto de herramientas integradas para ayudar al lingüista y, especialmente, favorecer el control de calidad mediante procedimientos automáticos para evaluar cada recurso, medir la regresión, verificar cierta armonía y coherencia entre los diferentes idiomas, gestión de las particularidades de los recursos multilingües.

<sup>1</sup> La sigla OAL corresponde a "Outils d'Aide aux Linguistes", es decir "Herramientas de Ayuda a los Lingüistas".

OAL ha sido diseñado como un frame-work cliente-servidor compuesto de:

- un servidor de recursos lingüísticos
- editores para la gestión de léxicos
- un agente inteligente de aspiración de contenido textual del Web
- una cadena de procesamiento para alimentar léxicos a partir de corpus brutos
- un *guesser* de paradigmas de flexión

La demo presenta el editor de léxicos morfosintácticos, sitúa esta herramienta en la arquitectura general de OAL y comenta los beneficios obtenidos por los lingüistas.

## 2 Editor de léxicos morfosintácticos

### 2.1 El formalismo SylLex

El formalismo SylLex permite modelar un léxico morfosintáctico. Este es concebido como un conjunto de lemas –simples y compuestos–, un conjunto de paradigmas de flexión y un conjunto de plantillas de paradigmas de flexión denominados *patterns*.

Un paradigma de flexión reagrupa un conjunto de reglas de flexión y genera las formas asociadas de un lema. La aplicación del paradigma "V/ar" al lema *comerciar*, p.ej., (cf. figura 1) genera todas las formas (flexiones) del verbo *comerciar*. Cada forma tiene asociada una etiqueta morfosintáctica que indica su categoría y atributos, como número o modo. Las formas se calculan a partir de fórmulas que

