Nicoletta Calzolari

Istituto di Linguistica Computazionale di Pisa

## <u>Lexicografía Computacional</u>

El procesamiento del lenguaje aplicado a corpus textuales a gramescala requiere la construcción de léxicos computacionales que contengan cientos de miles de entradas. Cada entrada debecontener de manera explícita información de todo tipo, incluida la semántica.

En este marco, se considera fundamental la reusabilidad de material lexicogràfico ya existente. El concepto de reusabilidad apareció por vez primera en el Workshop realizado en Grosseto (1986) y financiado por la CE, donde se planteó la necesidad do diseñar herramientas lingüísticas reusables, multifuncionales multilingües.

El concepto de reusabilidad debe entenderse en dos sentidos: a) para la extracción y uso de información léxica presente de manera explícita o implícita en fuentes de información léxic diversa (MRDs, Bases de Datos terminológicas, corpus textuales etc.) como ayuda para construir léxicos computacionales a gra escala y reusables en el sentido de b, y b) para construi léxicos computacionales destinados a diversos tipos de usuari (diferentes sitemas de PLN, lexicógrafos, lingüistas, usuarios n especializados).

Los léxicos computacionales actuales responden a uno de estos do tipos. La necesidad de disponer de recursos léxicos que responda a estas características se pone de manifiesto en el impulso dad por la CE a este tema, que se ha traducido en la puesta en march de diversos proyectos como son Acquilex (Esprit BRA), Eurotra 7 Multilex (Esprit) y Genelex (Eureka) y la subvención de otro como TEI (Text Encoding Initiatives) y Survey on Linguisti Resources.

El proyecto Esprit Acquilex puede considerarse como un prototip del primer grupo, mientras que otros proyectos, como Eurotra-7 como un prototipo del segundo.

## El proyecto Acquilex

El objetivo de la investigación en el marco de Acquilex se centren el desarrollo de técnicas y metodologías para extrae información, tanto explícita como implícita, de tipo sintáctico semántico, a partir de diccionarios en soporte magnético (MRDs para construir el componente léxico de sistemas de PLN.

Esta metodología se basa en la posibilidad de extraer de maner

automática la información que los MRDs contienen de manera implícita. Así, el diccionario se considera como una fuente de conocimiento básico y los objetivos fundamentales son la formalización de este conocimiento básico general en forma de conceptos y relaciones semánticas.

Entre los temas de investigación conectados con el objetivo diltimo de adquisición de conocimiento cabe destacar: a- el diseño de programas para la extracción de los superordinados a partir de las definiciones, su desambiguación y la construcción de taxonomías; b- el diseño de analizadores para el análisis de las definiciones con el objetivo de extraer la información semántica implícita; c- el estudio de como representar la información semántica extraída (i.e.: en forma de conceptos, atributos, (relaciones entre conceptos); d- el estudio de como relacionar las taxonomías y la representación conceptual; e- el diseño (implementación de software bàsico para la creación, acceso procesamiento de bases de datos léxicas y una base de conocimiento léxico.

Las taxonomías y las templetas (estructuras de rasgos en que si representa la información semántica) desarrolladas en el marco di Acquilex constituyen un primer grado de normalización estandarización en la representación de la semántica y de conocimiento general a partir de diccionarios (unos diez el total) y de lenguas diversas (cuatro).

En el proyecto Acquilex colaboran la universidades de Cambridge Dublin, Amsterdam, la Universidad Politéctica de Cataluña y e Instituto de lingüística Computacional de Pisa. Se trabaja sobr diccionarios monolingües del inglés, holandés, italiano castellano y sobre diccionarios bilingües inglés-holandés inglés-italiano.

Como complemento indispensable de los trabajos lexicográfico debemos citar las bases de datos textuales. En este marco cab destacar el programa TEI (Text Encoding Initiative) consistent en una iniciativa para definir criterios generales d codificación e intercambio de textos en soporte magnético. E trata de un proyecto internacional de cuatro años de duración qu tiene como objetivos: a)— especificar un formato uniforme par los textos en soporte magnético; b)— difundir las línea generales de codificación acordadas; y c)— documentar lo esquemas de documentación más importantes y desarrollar u metalenguaje en el que describirlos.