# Creating Tools for Morphological Analysis of Sumerian

## Valentin Tablan, Wim Peters, Diana Maynard, Hamish Cunningham

Department of Computer Science, University of Sheffield
Regent Court, 211 Portobello Street, S1 4DP, Sheffield, UK
{V.Tablan, W.Peters, D.Maynard, H.Cunningham}@dcs.shef.ac.uk

## The Sumerian Language

*The Sumerian language of ancient Sumer is a long-extinct language documented throughout the ancient Middle East, in particular in the south of modern Iraq, from at least the 4th millenium BC. It is arguably the first language for which we have written evidence, the rival candidate being ancient Egyptian. Sumerian was replaced by Akkadian as a spoken language around 2000 BC, but continued to be used as a sacred, ceremonial and scientific language in Mesopotamia until about 1 AD.*
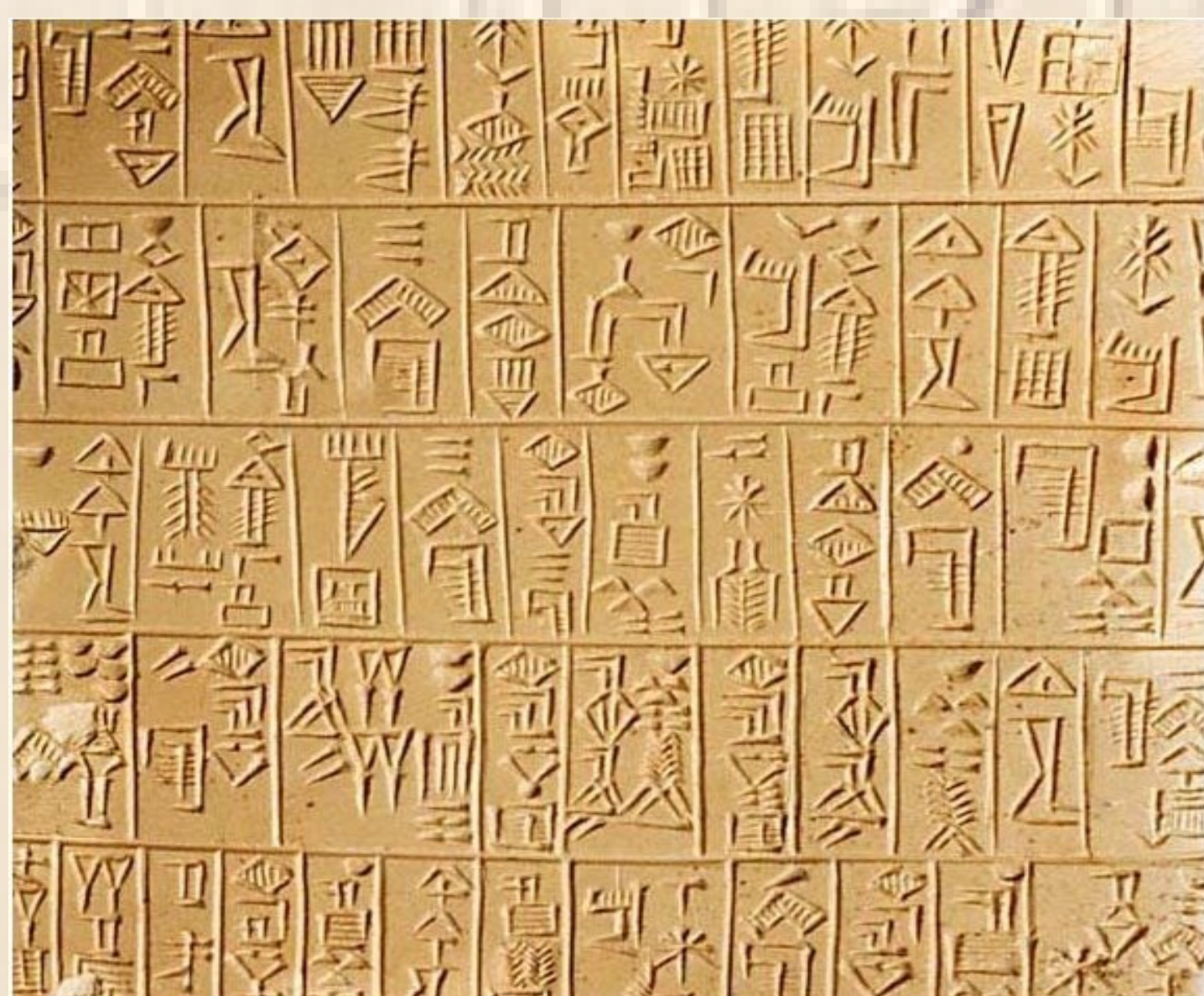
(Wikipedia, 2006)

## Sumerian Morphology

The Sumerian language is generally regarded as a language isolate in linguistics. No languages related to it have so far been convincingly identified, although many of its grammatical features are attested in other living languages outside of the Indo-European family to which English belongs.

In its orthographic form, Sumerian is encoded in cuneiform script, as depicted in this picture, which shows an example of 120 compartments of cuneiform script written by an expert scribe. Note that this example is much clearer and more beautiful than standard scripts because it describes the gifts from highly placed persons to a priestess.

Originally, cuneiform was logographic in nature, and a sign represented a content word (a thing or an action). It gradually developed into a combined system, where the same set of signs could be used to represent logograms and phonograms or syllabograms. In texts of the period we concentrate on, i.e. late third and early second millennium BC, logograms were used to write content words and the base (root) of a word, while phonograms were used to write bound morphemes and loan words.

In transliterated form, i.e. signs represented in the Roman alphabet with a few additions, these logograms and syllabograms are separated by a dash, as in `nam-lugal` (kingship), where the base root `lugal` (king) is combined with a derivational affix that changes the word into an abstract entity.

Another characteristic feature of Sumerian is the large number of homophones (words with the same sound structure but different meanings) - or perhaps pseudo-homophones, since there might have been differences in pronunciation (such as tone) that we do not know about. The different homophones (or, more precisely, the different cuneiform signs that denote them) are marked with different numbers by convention. For example: `du` = "to go", `du3` = "to build".

In terms of language typology, Sumerian is agglutinative. Word roots have grammatical elements glued on before or after them to build up complex grammatical forms. Many words (mostly verbs) consist of a root form (possibly reduplicated) and a chain of more or less clearly distinguishable and separable affixes or clitics. Nouns may have affix chains before as well as after the root. Overall, slightly less than 100 clitics have been postulated. Many of these clitics have allomorphs, depending on properties of the morphological context, such as progressive or regressive assimilation phenomena for vocals and consonants. On top of this, the null morpheme can occur as an allomorphic variant of a number of clitics. Overall, if we take all morphological rules into account, this leads to a very large number of possible interpretations.

Sumerian distinguishes the grammatical genders animate and inanimate, as do Polish, Russian, and some Native American languages, such as Navajo. There are also a large number of cases - nominative, ergative, genitive, dative, locative, comitative, equative ("as, like"), terminative ("to"), and ablative ("from").

## The ETCSL Project

The Electronic Text Corpus of Sumerian Literature - ETCSL , based at the University of Oxford, is the largest available corpus of literary Sumerian. It aims to make accessible on the web over 350 works composed during the late third and early second millennia BC. The corpus comprises Sumerian texts in transliteration, English prose translations and bibliographical information for each composition. The transliterations and the translations can be searched, browsed and read online using the tools of the website.

(http://www-etcsl.orient.ox.ac.uk/)

## Morphological Analysis Tools

The main aim of our work is to create a set of tools for performing automatic morphological analysis of Sumerian. This essentially entails identifying the part of speech for each word in the corpus (technically, this only involves nouns and verbs which are the only categories that are inflected), separating the lemma part from the clitics and assigning a morphological function to each of the clitics. In order to do this, we used the model of Sumerian morphology defined by a team of Sumerologists from the Unviersity of Oxford, which we then represented in a way that can be used for automated language processing.

The morphological model we used consists of noun and verb templates comprising a lemma plus a number of morphological slots that could be filled. The nouns have a lemma and up to six suffix slots while the verbs have up to twelve prefix slots, a lemma and two suffix slots. For each slot there is a known list of morphemes that can fill it and a set of restrictions encoding dependencies between the slots, such as agreement in gender. The lists of candidate slot fillers have non-null intersections – the same morpheme can appear in several lists, though usually with different Functions.

There are two main phases involved in our morphological analysis of transliterated Sumerian – a normalisation stage which deals with various surface phenomena which affect the way Sumerian words are written, such as reduplication or assimilation, and the actual morphological analysis which identifies parts of speech and assigns functions to the various morphemes.

Performance evaluation for the system in its current state, shows an F-Measure of 69% for noun identification, 66% for recognising verbs and 61% for morphological analysis of nouns.

## Corpus Search Tools

The linguistic analysis tools described above are complemented by the development of a tool for advanced search and visualisation of linguistic information, ANNIC (ANNotations In Context). This provides an alternative method of searching the textual data in the corpus, by identifying patterns in the corpus that are defined both in terms of the textual information (i.e. the actual content) and of metadata (i.e. linguistic annotation and XML/TEI markup).

ANNIC is based on Jakarta Lucene, but extends the model to allow users to query on annotations and their features by providing patterns that are similar to the JAPE rules in GATE (http://gate.ac.uk).