# A Data-driven Approach to the Mental Lexicon: Two Studies on Chinese Corpus Linguistics

Chu-Ren Huang, Kathleen Ahrens, and Keh-jian Chen

In this paper, we attempt to show i) that corpora offer real instances of language use (production) in a non-controlled environment, ii) that corpora constitute of a large sampling of the real input to linguistic perception, and iii) that corpora extracted from mass media represent the shared linguistic information of the language-speaking community.

Corpus-based studies are studies of linguistic theories based on linguistic objects (instead of on non-linguistic acts like naming, picture pointing, story-telling, or making decisions on yes-no questions.) We use two corpus-based studies to show that they can complement the traditional psychology-oriented studies based on controlled experiments. The two studies shed important light on the psychological reality of the notion of a word in the mental lexicon.

Our first study examines the definition of compounds based on M.I. (mutual information) values extracted from a corpus. We show that this empirically based definition of compounds easily resolves the previous controversies involving intuitive judgements (e.g. Bates et al. 1992 and 1993, and Zhou et al. 1993).

The second study involves the complex cognitive process of *suo1xie3* (abbreviation) and a simple statistical model. We show that while a rule-based model can only capture incomplete aspects of Chinese abbreviation, corpus-based statistical values nicely reflect their status in the mental lexicon.

In conclusion, we argue that corpora reflect shared uses of language and are efficient tools for establishing baseline facts in (psycho-/neuro-)linguistic research.

?    ?

?        ?