

Analogical Reasoning on Chinese Morphological and Semantic Relations

Shen Li ♠ Zhe Zhao Renfen Hu ♠ Wensi Li Tao Liu Xiaoyong Du
 ♠ {shen, irishu}@mail.bnu.edu.cn Beijing Normal University & Renmin University of China

GitHub: Chinese-Word-Vectors
 Over 2,000



Scan me

Introduction

Given the word representations, analogy questions can be automatically solved via vector computation:

$apples - apple + car \approx cars$ (morphological)
 $king - man + woman \approx queen$ (semantic)

It is well known that linguistic regularities vary a lot among different languages. For example, Chinese is a typical **analytic language** which lacks inflection.

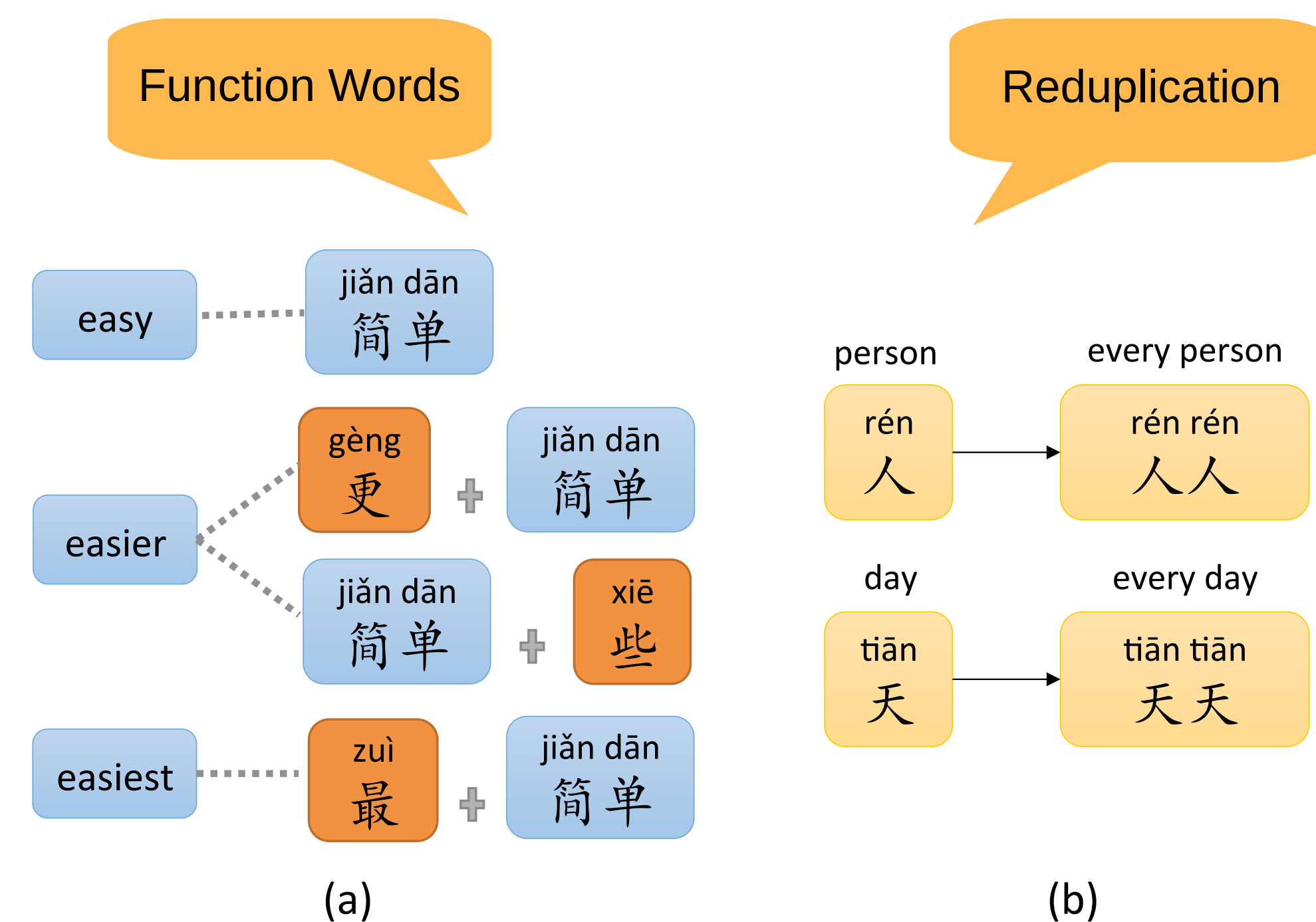
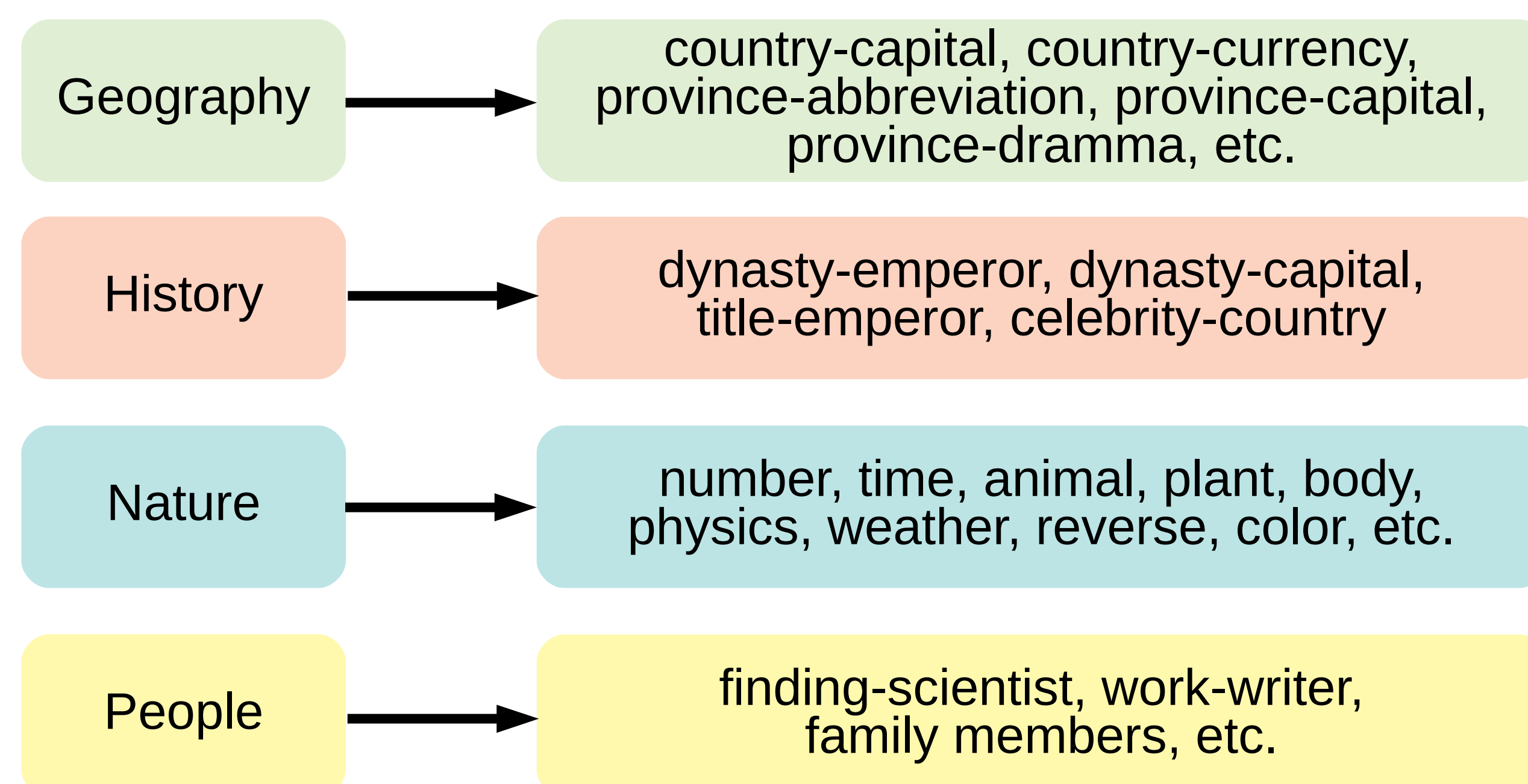


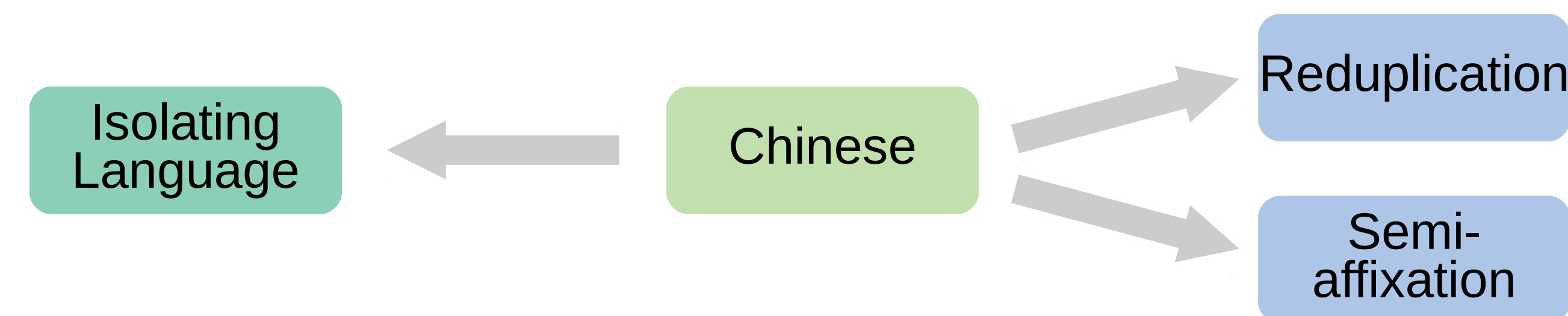
Figure 1: Examples of Chinese lexical knowledge.

Semantic Relations

We present **28 semantic relations** in 4 aspects.



Morphological Relations



• Reduplication

Reduplication means a morpheme is repeated to form a new word, which is semantically and/or syntactically distinct from the original morpheme.

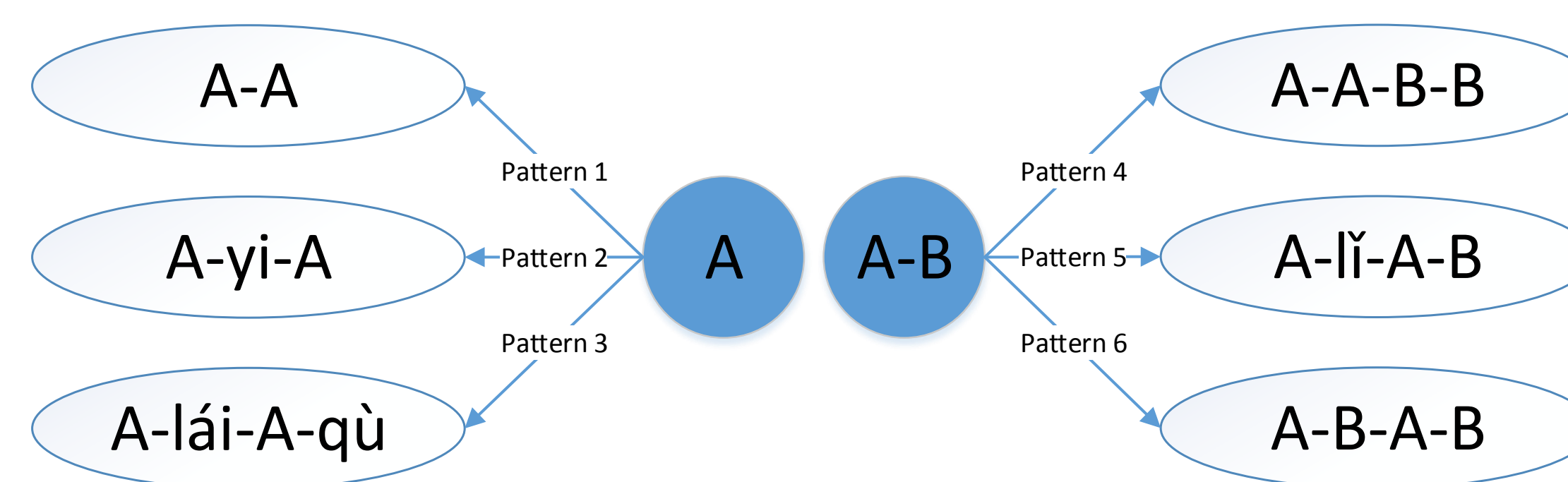


Figure 2: Reduplication patterns of A and A-B. (A and B are distinct morphemes.)

Taking $A \rightarrow AA$ as an example:

bà (dad) \rightarrow bà-bà (dad)
tiān (day) \rightarrow tiān-tiān (everyday)
shuō (say) \rightarrow shuō-shuo (say a little)
kàn (look) \rightarrow kàn-kàn (have a brief look)
dà (big) \rightarrow dà-dà (very big; greatly)
shēn (deep) \rightarrow shēn-shēn (deeply)

• Semi-affixation

Since Chinese is a typical isolating language that has few affixes, we describe the similar morphological process with semi-affixation suggested by Liu et al. 2001. To model the semi-affixation process, we uncover **21 semi-prefixes** and **41 semi-suffixes**.

Taking “dì-” and “-zi” as examples:

yī (one) \rightarrow dì-yī (first)
èr (two) \rightarrow dì-èr (second)
pàng (fat) \rightarrow pàng-zi (a fat man)
shòu (thin) \rightarrow shòu-zi (a thin man)

Pre-trained Embeddings

This project provides **100+ Chinese Word Vectors (embeddings)** trained with different representations (**dense and sparse**), context features (**word, ngram, character, and more**), and corpora.

Corpus	Size	Feature	Co-occurrence Type
Baidu Encyclopedia 百度百科	4.1G	Word	Word \rightarrow Word
Wikipedia_zh 中文维基百科	1.3G	Ngram	Word \rightarrow Ngram (1-2)
People's Daily News 人民日报	3.9G		Ngram (1-2) \rightarrow Ngram (1-2)
Sogou News 搜狗新闻	3.7G	Character	Word \rightarrow Character (1)
Financial News 金融新闻	6.2G		Word \rightarrow Character (1-2)
Zhihu_QA 知乎问答	2.1G		Word \rightarrow Character (1-4)
Weibo 微博	0.73G	Radical	Radical
Literature 文学作品	0.93G	Position	Word \rightarrow Word (left/right)
Mixed-large 综合	22.6G	Global	Word \rightarrow Word (distance)
Complete Library in Four Sections 四库全书	1.5G		Word \rightarrow Text
		Syntactic Feature	Word \rightarrow POS
			Word \rightarrow Dependency

Table 1: Corpus.

Table 2: Various Co-occurrence Information.

*All text data are preprocessed by removing HTML and XML tags. Only the plain text are kept and HanLP(v_1.5.3) is used for word segmentation.

CA8 Dataset

CA8 contains **17813** analogy questions and covers comprehensive morphological and semantic relations. CA8-morphological (**CA8-Mor**) contains **10177** morphological questions based on two types of relations: reduplication and semi-affixation. CA8-semantic (**CA8-Sem**) contains **7636** semantic questions divided into 4 categories and 28 sub-categories.