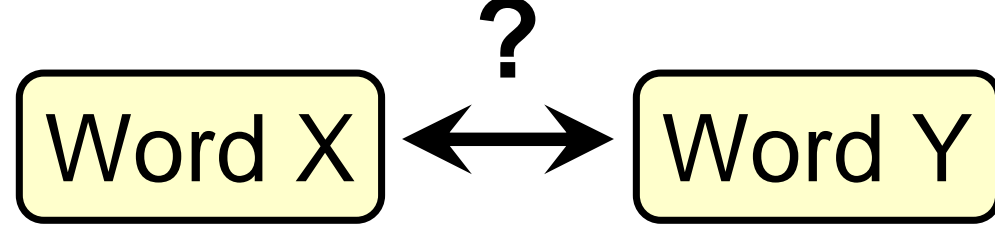
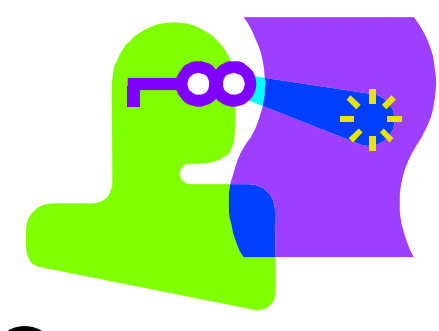


Statistical Analysis for Thesaurus Construction using an Encyclopedic Corpus

Yasunori OHISHI¹, Katunobu ITOU², Kazuya TAKEDA¹, Atsushi FUJII³
¹Nagoya University, Japan, ²Hosei University, Japan, ³University of Tsukuba, Japan

MOTIVATION

- Thesaurus
 - Expanding the range of possible terms in information retrieval
 - Example-based machine translation system
 - Most technical terminology is not included**
- Automatic thesaurus construction
 - Determining whether a word pairs has a hierarchical relation or a disrelation



PREVIOUS WORK

- To extract hyponyms, synonyms, and hypernyms,
 - Sentences that have specific syntactic patterns ("a part of" "is a" "such as") (Marti, 1992; Tsurumaru, 1991)
 - Descriptions in a dictionary (Suzuki, 2003)
 - Specific document structure (Shinzato, 2004) are used
- It is difficult to cover such an enormous vocabulary
 - Less frequent words are not expected in the desired expressions

THE CYCLONE CORPUS

- Collection procedure (Fujii, 2005)
 - Searching the Web for pages including a target term
 - Analyzing the layout of the pages and identifying paragraphs that potentially describe the target term
 - Classifying multiple paragraphs into predefined domains



Headword (Target term)

- Description 1
- Description 2

New terms and rarer terms

Many descriptions for each headword

- Semantic relations between a headword and terms in its description

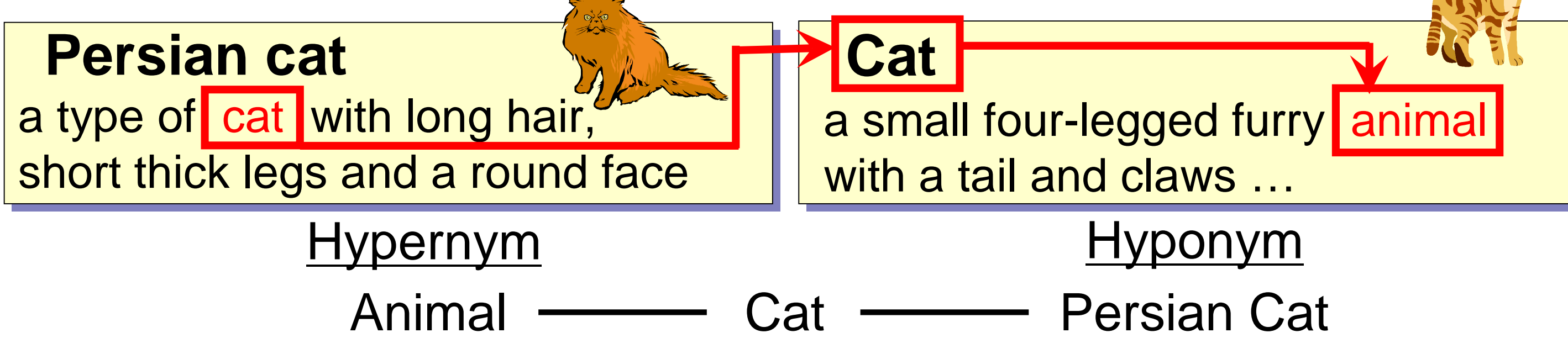
➔ Extracting hierarchical relations effectively

HIERARCHICAL RELATION ESTIMATION METHOD

- Descriptions of words have "directionality"
 - Lion**: a large wild animal of the cat family with yellowish brown fur ...
 - Animal**: a living creature such as a dog or cat that is generally distinguished from plants ...

A collection of descriptions for a headword shares common hypernyms but may not share hyponyms

- Target function $H(X|Y) = C(X|Y) - C(Y|X)$
- $C(X|Y)$: Probability of X , given the descriptions of Y
- Semantic expansion technique for descriptions



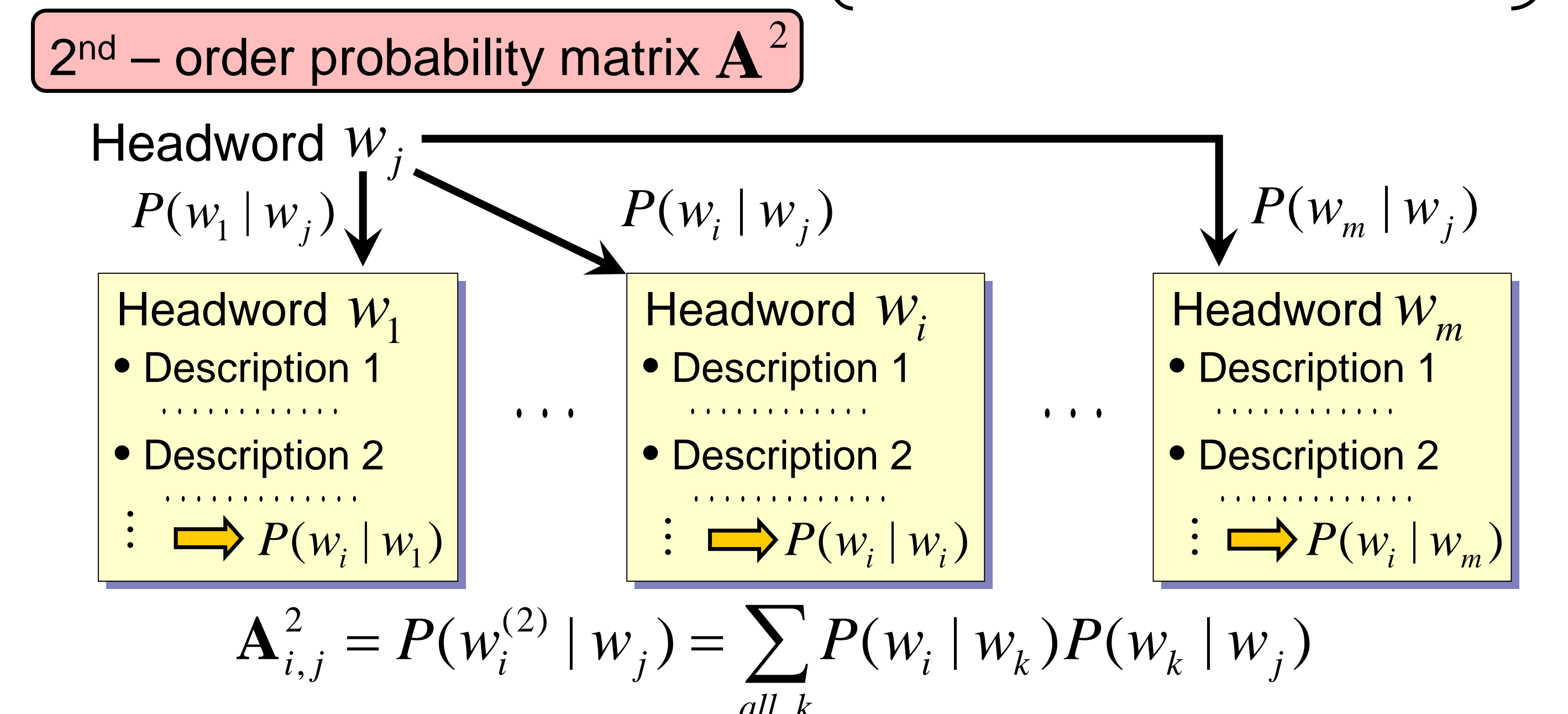
1st - order probability matrix \mathbf{A}

$$\mathbf{A}_{i,j} = P(w_i | w_j)$$

$$i, j = 1, \dots, m$$

Relative frequency: $P(w_i | w_j) = \frac{F(w_i | w_j)}{\sum_{k=1}^m F(w_k | w_j)}$

m is the number of headwords



n^{th} - order probability matrix is defined as a square matrix \mathbf{A}^n

Expanded probability matrix

$$\mathbf{C} = \sum_{i=1}^N [\alpha_i \mathbf{A}^i] \quad [N \text{ is the maximum order of expansion}]$$

Identifying the relation of a word pair (w_i, w_j) by calculating $H(w_i | w_j)$

EVALUATION

- Word pairs in the computer-related domain (2074 words)

Correctness of description	Average number of descriptions per headword	Test sets	
		Hierarchical relation data	Disrelation data
A	6.6	136 pairs	301 pairs
A+B	10.4	168 pairs	366 pairs
ALL	80.7	206 pairs	497 pairs

(A: Correct, B: Partially correct)

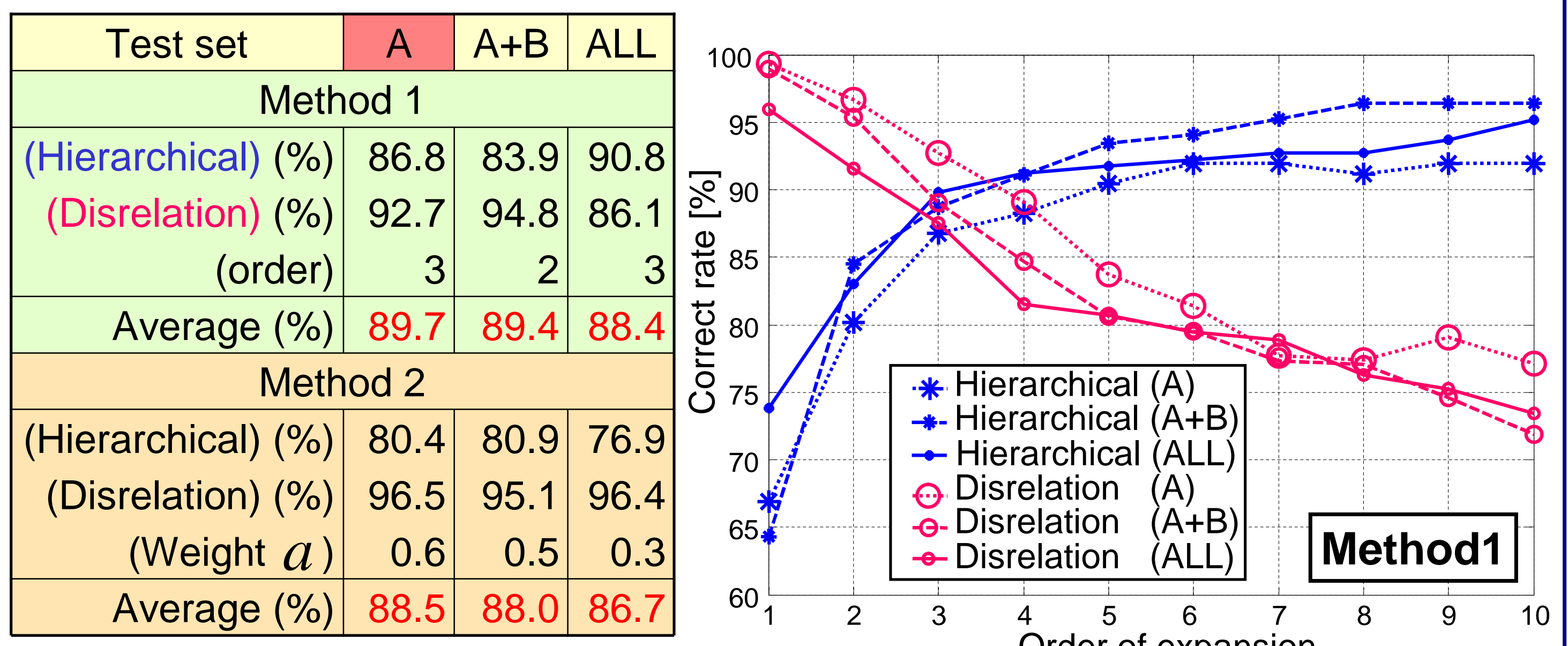
Example

Hierarchical relation data		Disrelation data	
Hypernym	Hyponym		
ISO	MPEG	Bit error rate	MS-DOS
Operations research	Linear programming	Polish notation	10BASE-T
OS	UNIX	Compiler	Zoned decimal
Protocol	TCP/IP	Hexadecimal number	CATV
Memory	DRAM	Dual system	Graph theory

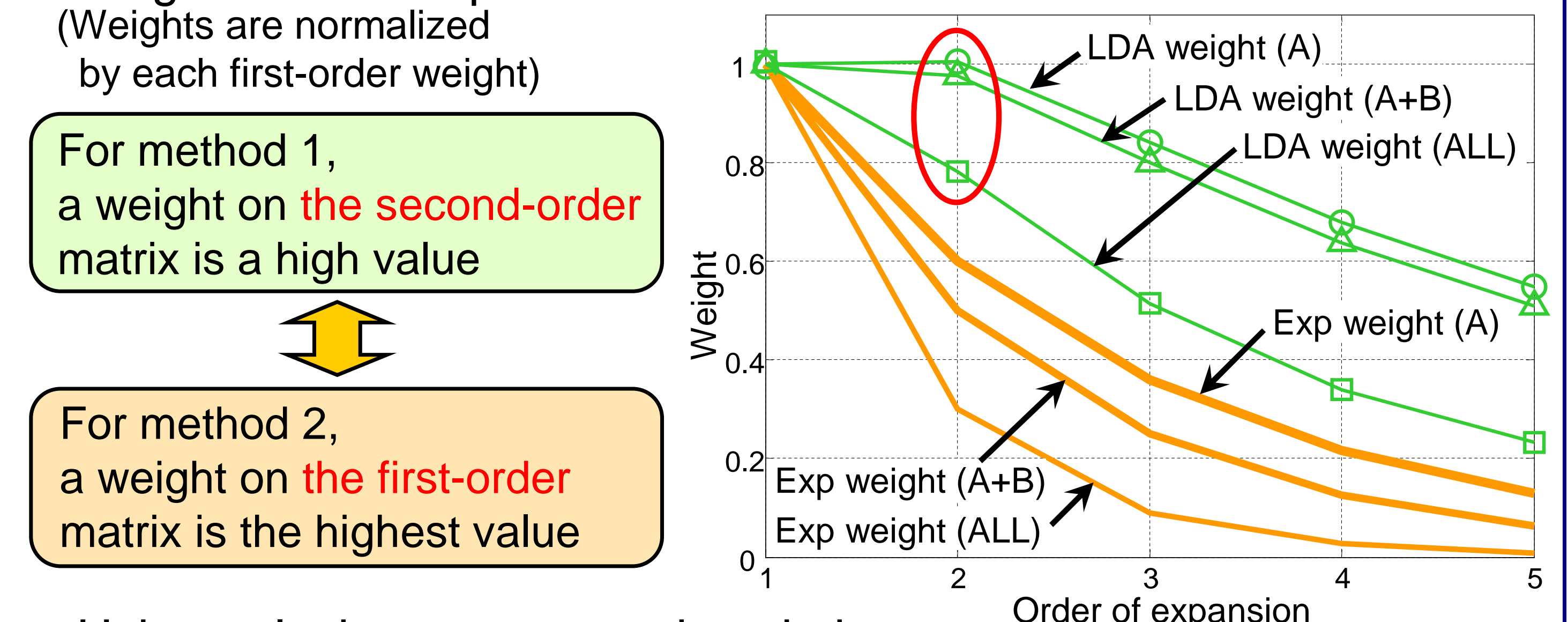
- To determine weights $\alpha_i, (i = 1, \dots, N)$
- Method1** Estimation using linear discrimination analysis (proposed method)
- Method2** Exponential weight method (Suzuki, 2003)

$$C = \lim_{k \rightarrow \infty} b(aA + a^2A^2 + \dots + a^kA^k) \quad 0 < a < 1, b = \frac{1-a}{a}$$

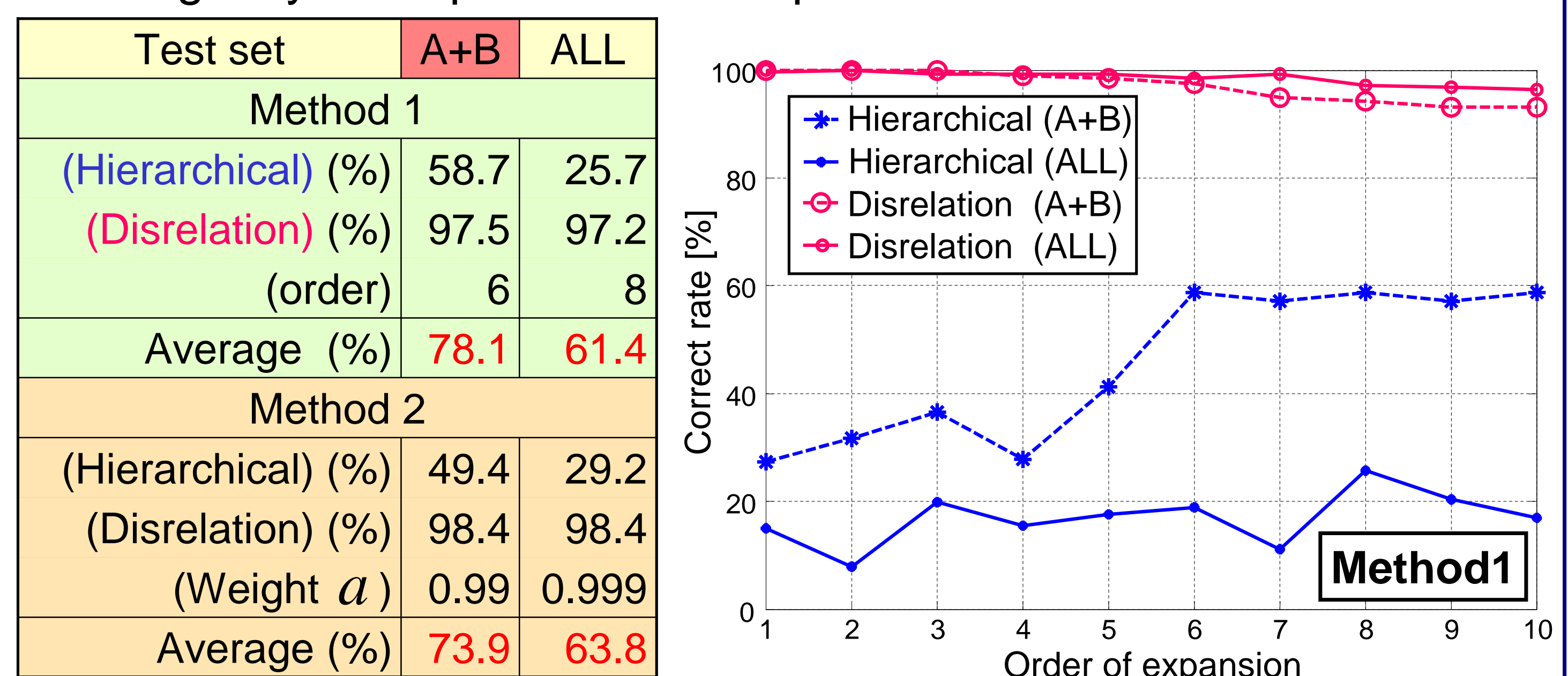
- Hierarchical relation or disrelation?



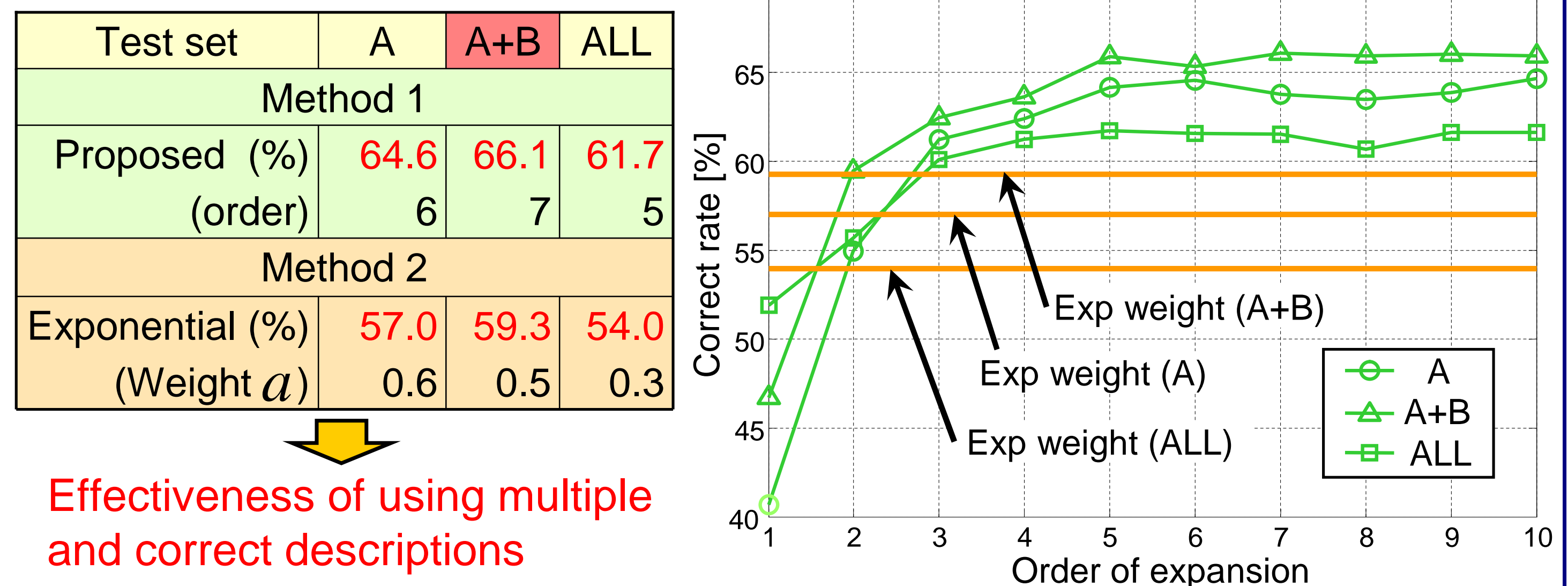
- Weights for the expanded matrices



- Using only the top correct description



- Hierarchical discrimination



CONCLUSION AND FUTURE WORK

- Conclusion**
 - Discrimination for the hierarchical relation of a word pair using an encyclopedic corpus called the Cyclone corpus
 - In order not to miss an indirect relationship, a semantic expansion technique for descriptions is used
 - The proposed method is able to detect 66.1% of relations
- Future work**
 - Discrimination between hierarchical and synonymous relation