

< IWP 2005, Oct. 14th, 2005 >

## A Class-oriented Approach to Building a Paraphrase Corpus

Atsushi FUJITA<sup>(1)</sup>, Kentaro INUI<sup>(2)</sup>

<sup>(1)</sup> Kyoto University

<sup>(2)</sup> Nara Institute of Science and Technology

## Requirements for handling paraphrases

### Transformation rules / patterns

- Handcrafting  $X_{verb} \Leftrightarrow S_d(X) + Oper_f(S_d(X))$   
[Iordansjaka et al., 1991] [Dras, 1999] [Sato et al., 1999]  
[Kondo et al., 1999] [Kondo et al., 2001] [Iida et al., 2001] etc.

- Automatic acquisition  
[Barzilay et al., 2001] [Lin et al., 2001] [Shinyama et al., 2002]  
[Shimohata et al., 2002] [Pang et al., 2003] etc.

burst into tears  $\Leftrightarrow$  cry    X finds a solution to Y  $\Leftrightarrow$  X solves Y

### Paraphrase corpus (collection of paraphrase examples)

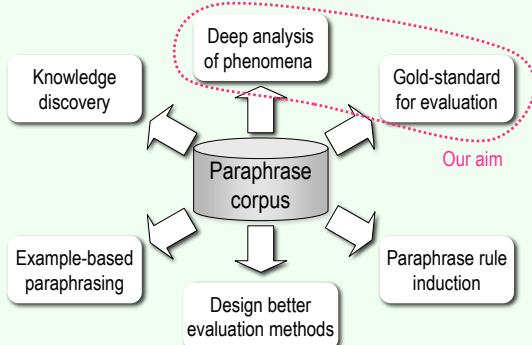
- Few freely available resources [Dolan et al., 2004]

The leading indicators measure the economy...  
The leading index measures the economy....

2

## Purposes of paraphrase corpus

### Beneficial to activate the research field



3

## Outline

- Background
- Issues and our class-oriented approach
- Semi-automatic example collection
- Preliminary trials
  - Specification
  - Discussion
- Conclusion

4

## Building paraphrase corpus

### Issues

- to consider: variety, source, organization
- to maximize: coverage, reliability, cost-efficiency

### Previous work

#### Manual production

[Shirai et al., 2001]  
[Kinjo et al., 2003]  
[Shimohata et al., 2004]

#### Automatic acquisition

[Barzilay et al., 2003]  
[Shinyama et al., 2003]  
[Dolan et al., 2004]

Reliability  Cost-efficiency

- Coverage is not ensured
- No focus on sorts/variety of paraphrases

5

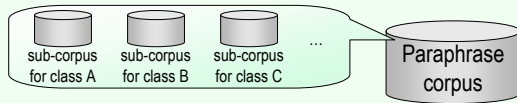
## Variety of paraphrases



6

## A class-oriented approach

### Separately collect examples for each class



### Semi-automatic example collection

- Automatic generation + human judgment
- Step 1:** Define a paraphrase class based on morpho-syntactic transformation patterns
- Step 2:** Collect all candidates using a paraphrase engine
- Step 3:** Judge candidate paraphrases in hand

7

## Aim of this study

### Confirm the feasibility of the method through practice

- Given
  - A paraphrase class
  - A text collection
- Collect paraphrase examples belonging to the class
  - At a **minimal human labor cost**
  - As **exhaustively** as possible from the text collection
  - As **reliable** as humanly possible

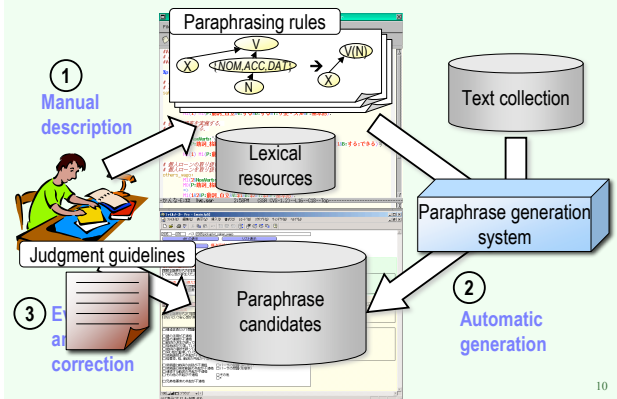
8

## Outline

- Overview
- Issues and our class-oriented approach
- Semi-automatic example collection
  - Specification
  - Discussion
- Preliminary trials
- Conclusion

9

## Semi-automatic example collection



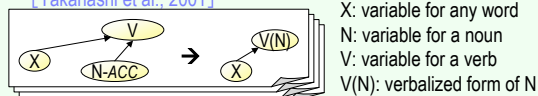
10

## Step 1: Pattern description

### Morpho-syntactic paraphrasing patterns

- Pairs of dependency trees
- Implemented on a paraphrase generation system

[Takahashi et al., 2001]



映画が 彼に 感動を 与える  
 film-NOM him-DAT impression-ACC to give-ACTIVE  
 (The film made an impression on him.)

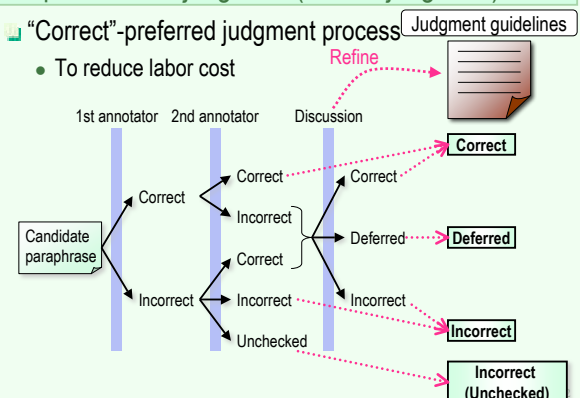
映画が 彼を 感動させる  
 film-NOM him-ACC to be impressed-CAUSATIVE  
 (The film impressed him.)

11

## Step 3: Manual judgment (mutual judgment)

### "Correct"-preferred judgment process

- To reduce labor cost



## Step 3: Manual judgment (I/F)

The screenshot shows a web-based interface for manual judgment. Annotations include:

- (a) source sentence: Points to the 'Given' text area.
- (b) automatically generated paraphrase: Points to the 'Generated' text area.
- (c) second opinion (correct / incorrect): Points to the 'Obligatory' checkbox.
- (c) annotator's judge (correct / incorrect): Points to the 'Optional' checkbox.
- (d) error tags: Points to the 'Error tags' input field.
- (e) confirmed (revised) paraphrase: Points to the 'Confirmed' checkbox.
- (f) free comments: Points to the 'Free comments' text area.

13

## Outline

1. Overview
2. Issues and our class-oriented approach
3. Semi-automatic example collection
4. Preliminary trials
  1. Specification
  2. Discussion
5. Conclusion

14

## Target classes

### Paraphrases of light-verb constructions (LVC)

映画<sup>が</sup> 彼<sup>に</sup> 感動<sup>を</sup> 与<sup>え</sup>る  
 film-NOM him-DAT impression-ACC to give-ACTIVE  
 (The film made an impression on him.)

映画<sup>が</sup> 彼<sup>を</sup> 感動<sup>さ</sup>せる  
 film-NOM him-ACC to be impressed-CAUSATIVE  
 (The film impressed him.)

### Transitivity alternation (TransAlt)

そよ風<sup>が</sup> 木々<sup>を</sup> 揺<sup>ら</sup>す  
 breeze-NOM tree-ACC to sway-Transitive  
 (The breeze sways the trees.)

木々<sup>が</sup> そよ風<sup>に</sup> 揺<sup>れ</sup>る  
 tree-NOM breeze-DAT to sway-Intransitive  
 (The trees sway in the breeze.)

15

## Resources

### LVC

- 4 paraphrasing patterns (e.g. (7) in paper)
- 20,155 pairs of <deverbal noun, verb>
  - 感動(impression), 感動する(to be impressed)
  - 誘い(invitation), 誘う(to invite)

### TransAlt

- 8 paraphrasing patterns (e.g. (10) in paper)
- 212 pairs of <intransitive verb, transitive verb>
  - 揺れる(to sway-Intransitive), 揺らす(to sway-Transitive)
  - 壊れる(to break-Intransitive), 壊す(to break-Transitive)

16

## Results of trials

	Paraphrase class	LVC	TransAlt
Step 1	# of paraphrasing patterns	4	8
	Size of dictionary	20,155	212
Step 2	# of source sentences	10,000	25,000
	# of generated candidates	2,566	985
Step 3	# of judged candidates	1,067	964
	# of incorrect candidates	520	503
	# of correct candidates	547	461
	# of paraphrase examples	591	484
	Working hours	118	169.5

Working hours: 2 annotators' working time for (1) Judgment, (2) Discussion, and (3) Re-judgment after refining guidelines

17

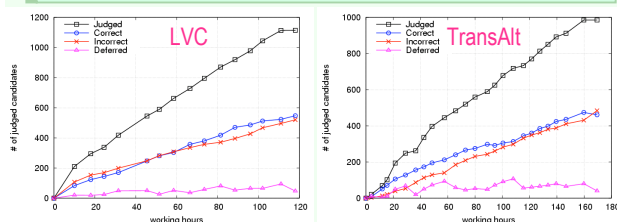
## Aim of this study (reminder)

### Confirm the feasibility of the method through practice

- Given
  - A paraphrase class
  - A text collection
- Collect paraphrase examples belonging to the class
  - At a minimal human labor cost
  - As exhaustively as possible from the text collection
  - As reliable as humanly possible

18

## Cost-efficiency



### Not so wasteful human labor cost

- 7.1 candidates / man-hour
- 3.7 paraphrase examples / man-hour
- TransAlt is 1.75 times more difficult than LVC due to test

19

## Exhaustiveness

### The initial resource is not necessarily optimal

- Paraphrasing patterns
- Derivation pairs

### How are they optimal?

- Estimated coverage: 77% ( $158/(158+48)$ )
  - 158 paraphrases for 750 excerpted sentences
  - Manual examination obtained another 48 paraphrases
    - 47 misses can be salvaged by resource enhancement
    - Errors of shallow parsers hurt only once
- Use of patterns is realistic approach
- Manual examination ensures coverage

20

## Reliability

### Strategy

- Classification bases on guideline & linguistic intuition
- Inter-annotator discussion refined judgment guidelines

### Agreement ratio increased (in case of LVC)

- 74% (3<sup>rd</sup> day) → 77% (6<sup>th</sup> day)  
→ 88% (9<sup>th</sup> day) → 93% (11<sup>th</sup> day)
- It's still not easy to explain "why this is correct / incorrect"

### Future plan

- Involve an expert to make sure of judgment guidelines
- Involve the 3<sup>rd</sup> annotator for judgment

21

## Conclusion

### Feasibility of a semi-automatic example collection

- Class-oriented example collection
- Employing a paraphrase generation system

### Promising results

- Reasonable human labor cost, but need reduction
- Moderately exhaustive at initial stage
- Typically reliable, but some marginal cases

### Paraphrase sub-corpora consist of

- LVC: 1067 candidates / 591 examples
- TransAlt: 964 candidates / 484 examples

22

## Future work

### Discussion on required expertise

- It is not easy to explain "why this is correct / incorrect"
- Involve an expert to make sure of judgment guidelines

### Build sub-corpora for other paraphrase classes

### Extrinsic evaluation through case studies

- Resultant provides both correct and incorrect examples
- Immediately available for analysis and system evaluation

### Publicly open the resource

- Paraphrase corpus, Lexical resources, Judgment guidelines

23