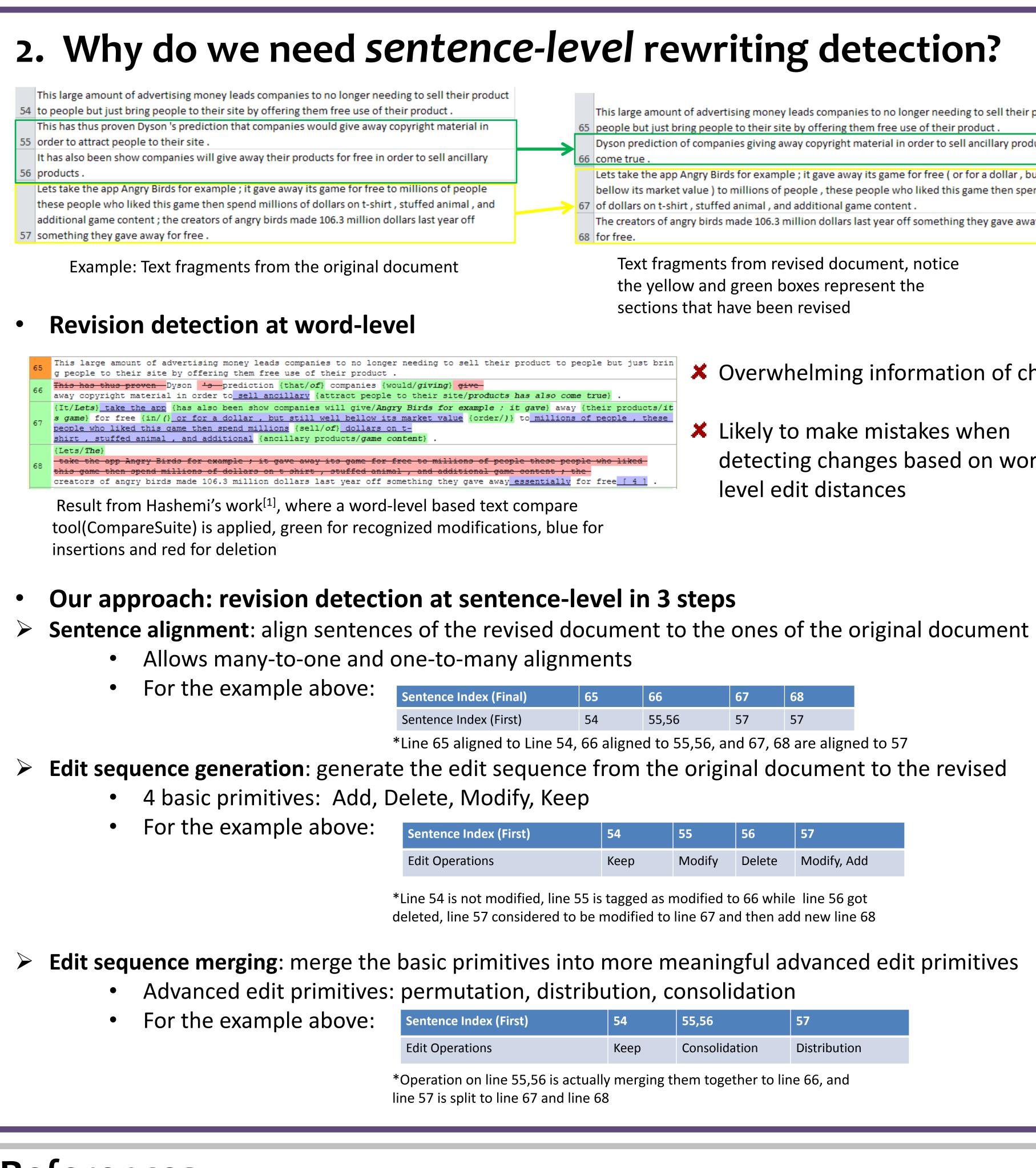
# **Sentence-level Rewriting Detection**

Fan Zhang<sup>1</sup>, Diane Litman<sup>1,2</sup> {zhangfan, litman}@cs.pitt.edu

# 1. Our goal

- Current goal: Help to have a better understanding of the rewriting process
- Ultimate goal: Provide automatic revision suggestions



## References

[1] Homa B. Hashemi and Christian D. Schunn. 2014. A tool for summarizing students' changes across drafts. In International Conference on Intelligent Tutoring Systems (ITS) [2] Cho K, Schunn C D. Scaffolded writing and rewriting in the discipline: A web-based reciprocal peer review system[J]. Computers & Education, 2007, 48(3): 409-426. [3] Rani Nelken and Studart M Shieber. 2006. Towards robust context-sensitive sentence alignment for monolingual corpora. In EACL [4] Saul B Needleman and Christian D Wunsch. 1970. A general method applicable to the search for similarities in the amino acid sequence of two proteins. Journal of molecular biology, 48(3): 443-453

# <sup>1</sup>Department of Computer Science, University of Pittsburgh | <sup>2</sup>Learning Research and Development Center, University of Pittsburgh

	el rewriting detection?
5	This large amount of advertising money leads companies to no longer needing to sell their product to people but just bring people to their site by offering them free use of their product .
6	Dyson prediction of companies giving away copyright material in order to sell ancillary products has also come true .
	Lets take the app Angry Birds for example ; it gave away its game for free (or for a dollar, but still well bellow its market value) to millions of people, these people who liked this game then spend millions of dollars on t-shirt, stuffed animal, and additional game content. The creators of angry birds made 106.3 million dollars last year off something they gave away essentially for free.
	Text fragments from revised document, notice the yellow and green boxes represent the sections that have been revised
op	De but just brin Overwhelming information of changes

**X** Likely to make mistakes when detecting changes based on wordlevel edit distances

65	66	67	68
54	55,56	57	57

\*Line 65 aligned to Line 54, 66 aligned to 55,56, and 67, 68 are aligned to 57

54	55	56	57
Кеер	Modify	Delete	Modify, Add

\*Line 54 is not modified, line 55 is tagged as modified to 66 while line 56 got deleted, line 57 considered to be modified to line 67 and then add new line 68

	54	55,56	57		
	Кеер	Consolidation	Distribution		
ually merging them together to line 66, and					

<ul> <li>J. Our</li> <li>Data p</li> <li>2 under</li> </ul>
Collecte     Corpus1     Corpus2
<ul> <li>Manua</li> <li>Sentena</li> <li>Edit sea</li> <li>Edit sea</li> </ul>
<ul> <li>Autom</li> <li>Sentend</li> <li>Method: ad</li> <li>Logistic</li> <li>Global</li> <li>Evaluat</li> </ul>
<ul> <li>Performane</li> <li>Baseline</li> <li>Group</li> <li>Cross validation</li> <li>Cross validation</li> <li>Train on corpus</li> <li>Train on corpus</li> </ul>
<ul> <li>*All the algorit</li> <li>Edit seq</li> <li>Method:</li> <li>Evaluatio</li> <li>Performa</li> <li>Baseline</li> </ul>
Baseline Rule-based met Rule-based on a * Rule-base app
still better than <b>Edit seq</b> Method: Rule Evaluation: a Performance
<ul> <li><b>4.</b> Fut</li> <li>Improve</li> <li>Replace</li> <li>distance</li> <li>Identify</li> </ul>

This research is supported by the Institute of Educational Sciences, U.S. Department of Education, through Grant R305A120370 to the University of Pittsburgh. The opinions expressed are those of the authors and do not necessarily represent the views of the Institute or the U.S. Department of Education



# **2. Our work**

#### reparation

rgraduate paper assignments from a "Social Implications of Computing Technology" course ed via a web-based peer review system<sup>[2]</sup>, each paper has two drafts

# of pairs	# of sentences in the first draft			Average words in one sentence in C2
11	761	791	22.5	22.7
10	645	733	24.7	24.5

### annotation

ce alignment: two annotators annotate on one paper, kappa: 0.794 quence generation: annotate edit sequence from the first draft quence merging: annotates "consolidation", "permutation" currently

### natic sentence-level revision detection in 3 steps

#### ce alignment

dapting Nelken's approach[3]

regression classifier using sentence similarity score (Word Overlap, TF\*IDF, Levenshtein Distance) alignment based on sentence order (Needleman-Wunsch<sup>[4]</sup>)

- tion: accuracy (percentage of sentences that are correctly aligned)

e: Hashemi's word-based approach (as in section 1), performance collected by manual inspection

	Levenshtein Distance	Word Overlap	TF*IDF	Baseline
n on corpus 1	0.9811	0.9863	0.9931	0.9427
n on corpus 2	0.9649	0.9593	0.9667	0.9011
and test on corpus 2	0.9727	0.9700	0.9727	0.9045
2 and test on corpus 1	0.9860	0.9886	0.9798	0.9589

thms achieved high accuracy, TF\*IDF achieves the best performance among all similarity metrics

### uence generation

Rule-based approach

on: Word error rate (WER), rate of segments to be modified to match with the correct sequence nce

e: Hashemi's word-based approach (as in section 1), performance collected by manual inspection

	Corpus1	Corpus2
	0.091	0.153
ethod	0.035	0.017
alignment results	0.067	0.025

proach achieved a much better performance of edit sequences comparing to baseline, the result is the baseline even applied on the automatic aligned results from the first step

#### juence merging

le-based approach, now only recognizes "Distribution" and "Consolidation" accuracy (percentage of the "Distribution" and "Consolidation" cases recognized) e: The 9 consolidation and 5 distribution cases are all successfully recognized (100% accuracy)

## ture work

the accuracy of current algorithm

rule-based approach used in edit sequence generation phase with approach based on edit , and then infer advanced edits based on the automatic generated sequence

more meaningful advanced rewriting operations

Conduct user study comparing the utility of sentence versus word-level rewriting detection

