

# Dependency Forest for Sentiment Analysis

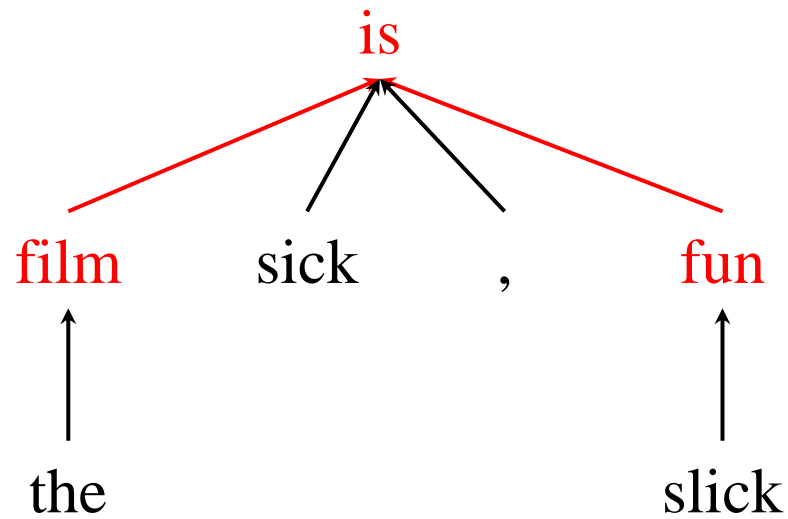
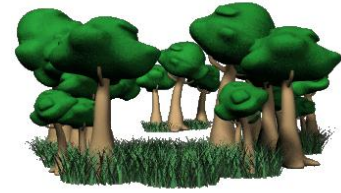
Zhaopeng Tu, Wenbin Jiang, Qun Liu, Shouxun Lin

*Key Laboratory of Intelligent Information Processing  
Institute of Computing Technology*



中国科学院  
INSTITUTE OF COMPUTING  
TECHNOLOGY

# Dependency Grammars in SA

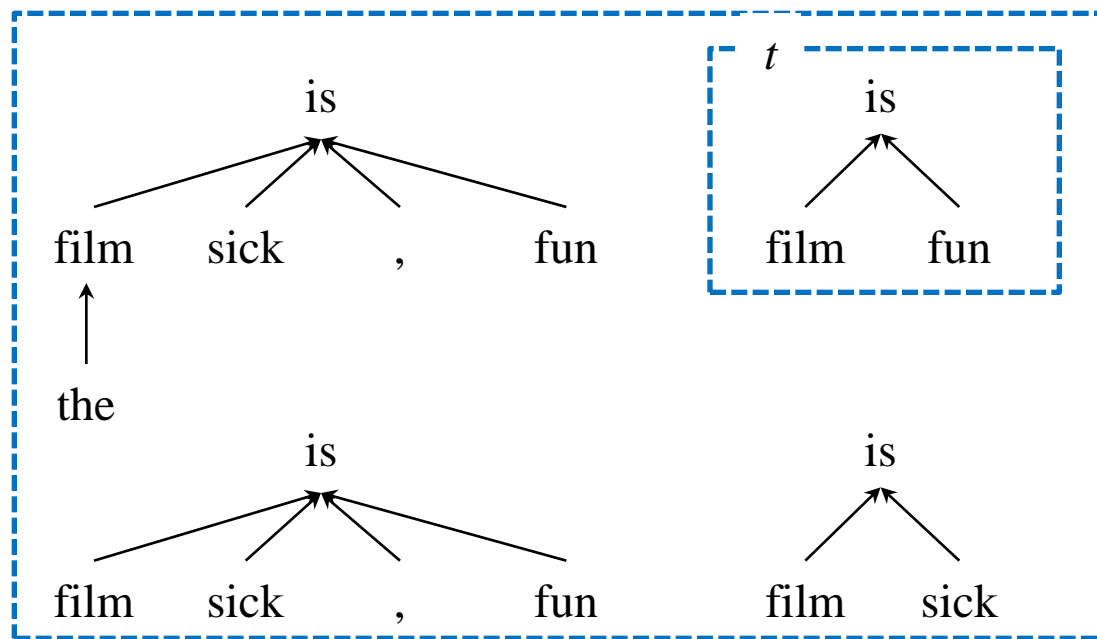
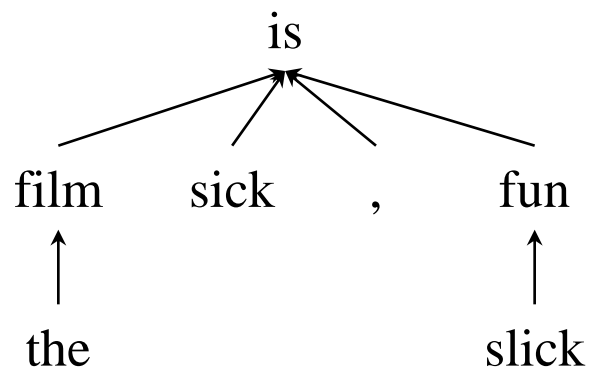


(Matsumoto et al., 2005; Joshi et al., 2009; Liu et al., 2009)

# Tree-based Features Extraction

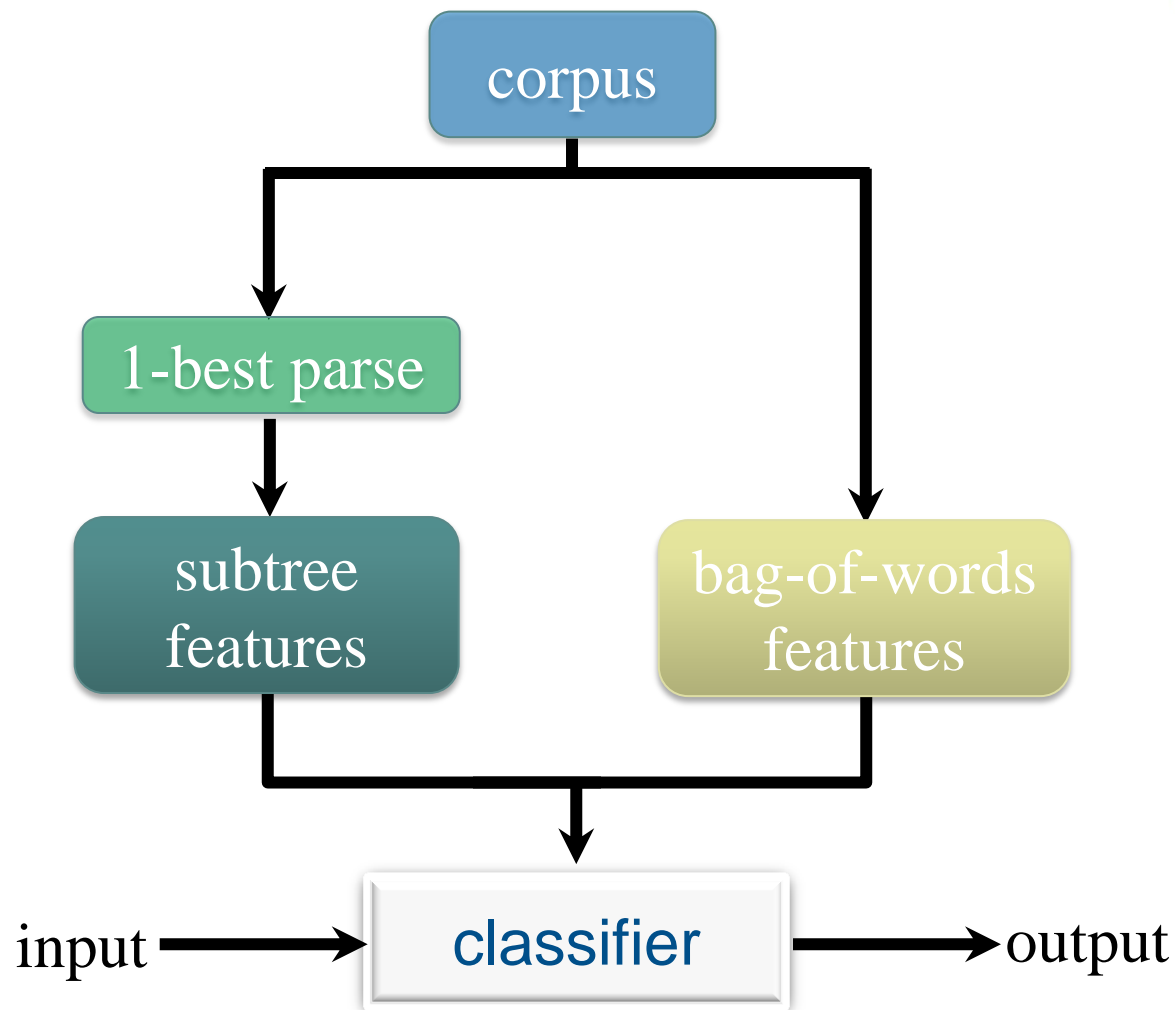


(Matsumoto et al., 2005)

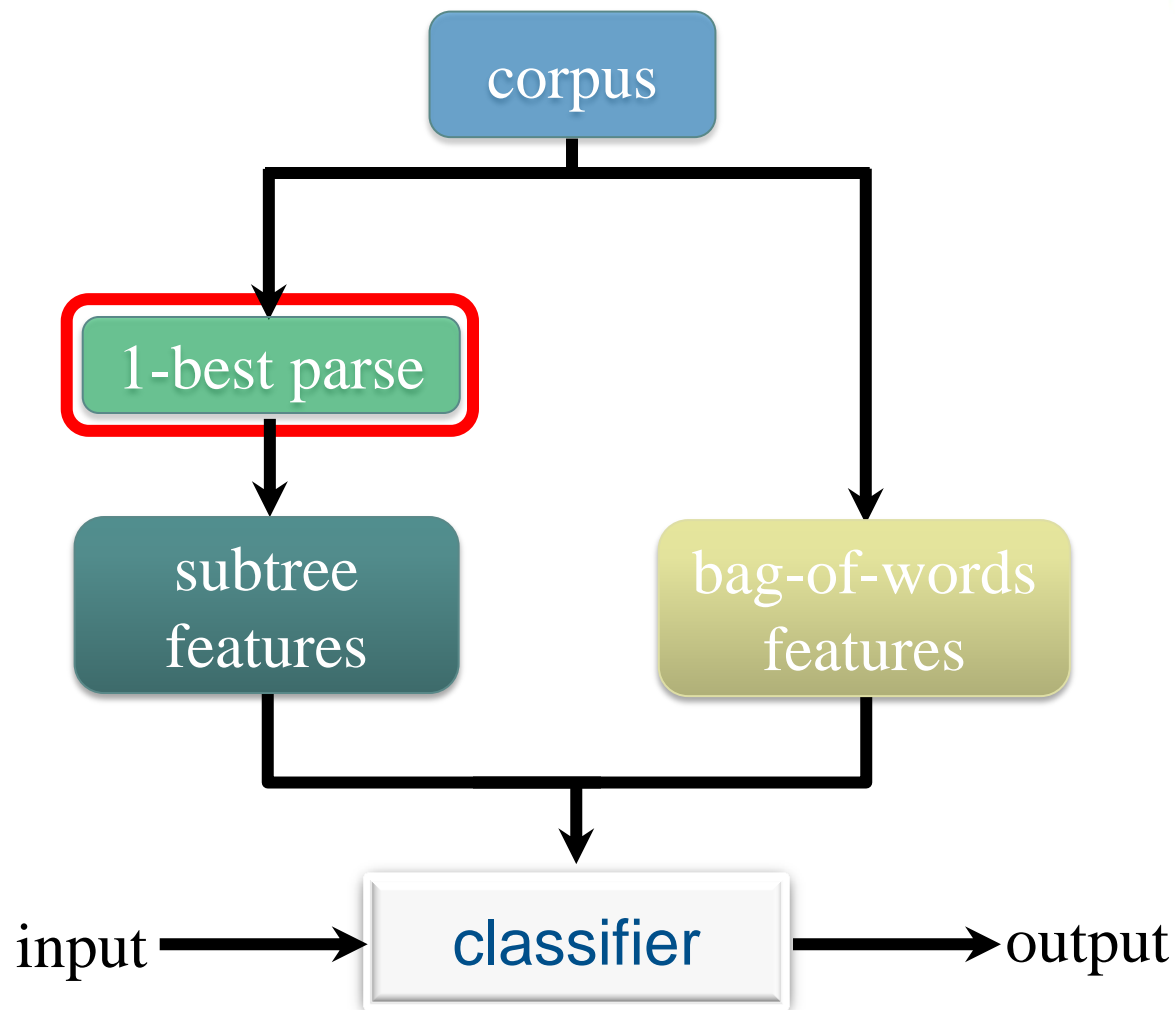


extracted subtrees

# Pipeline



# Pipeline



# Challenges



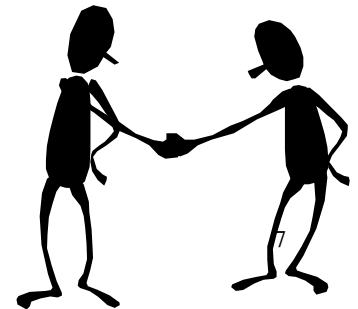
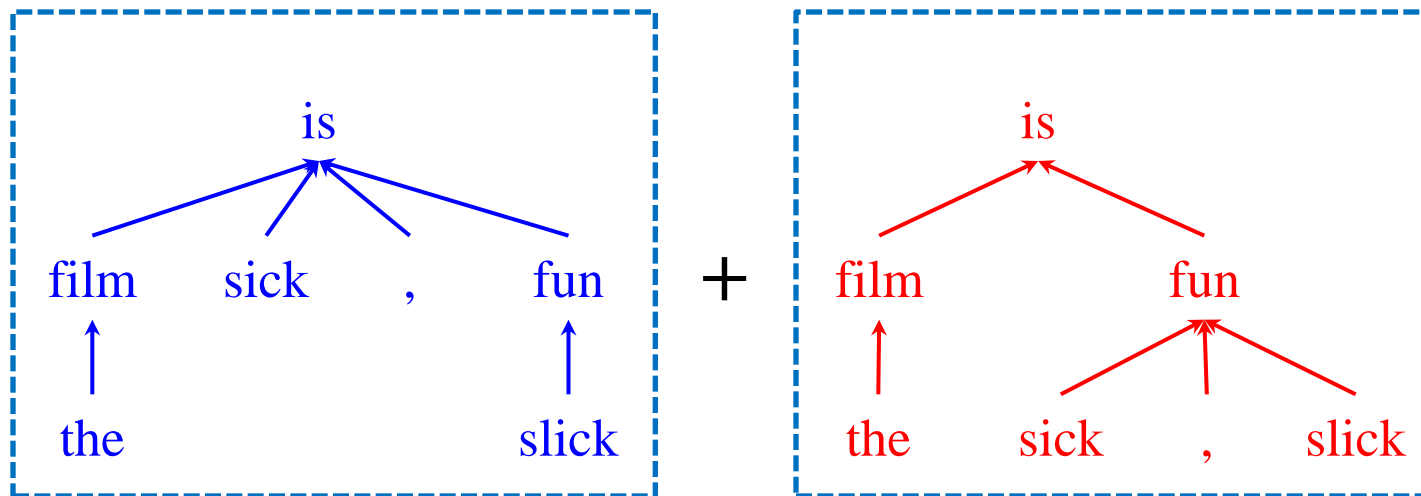
- Tree-based approach faces the major challenges:
  - vulnerable to parsing error



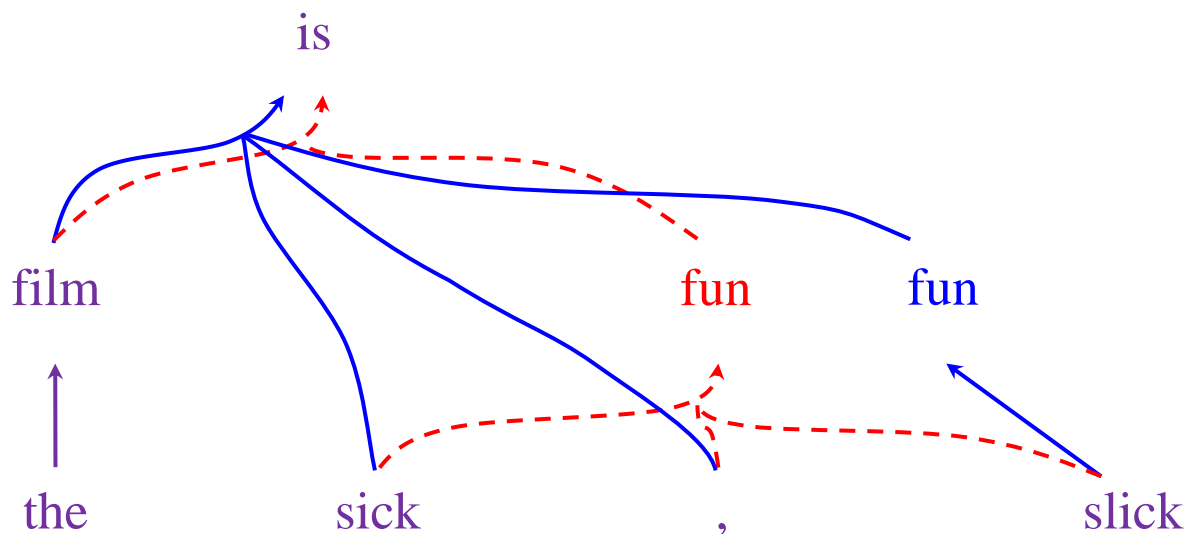
# Solution



- Dependency forest provides an elegant solution to this problem (Tu et al, 2010)

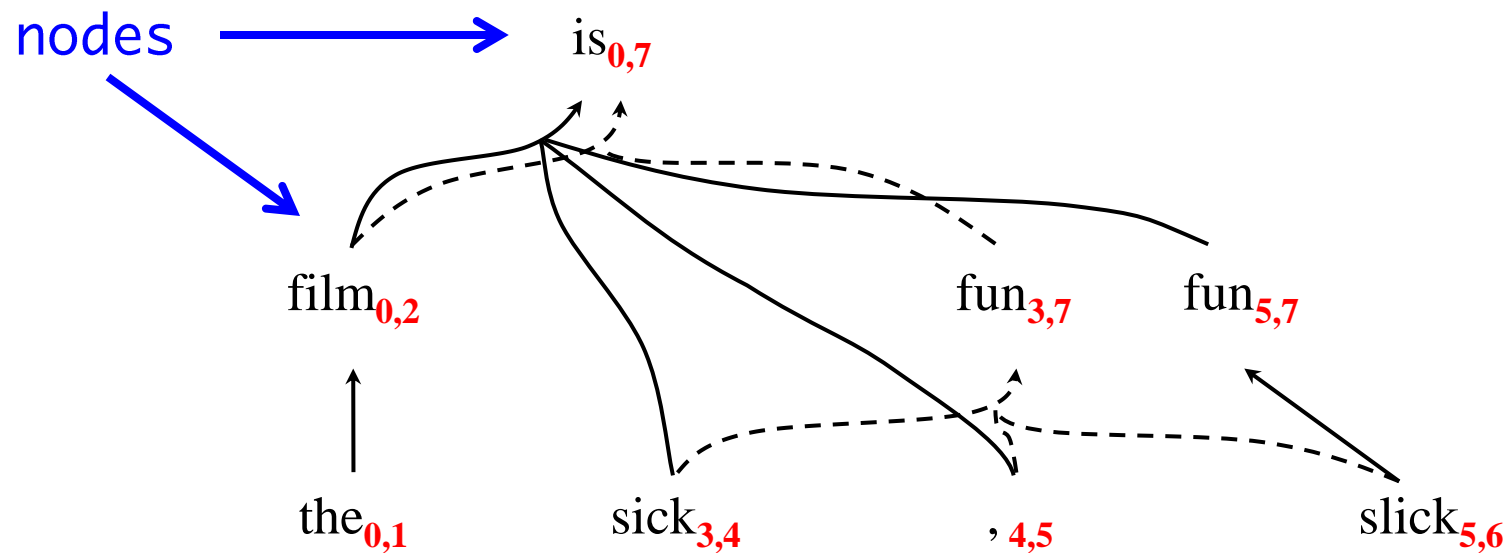


# Dependency Forest

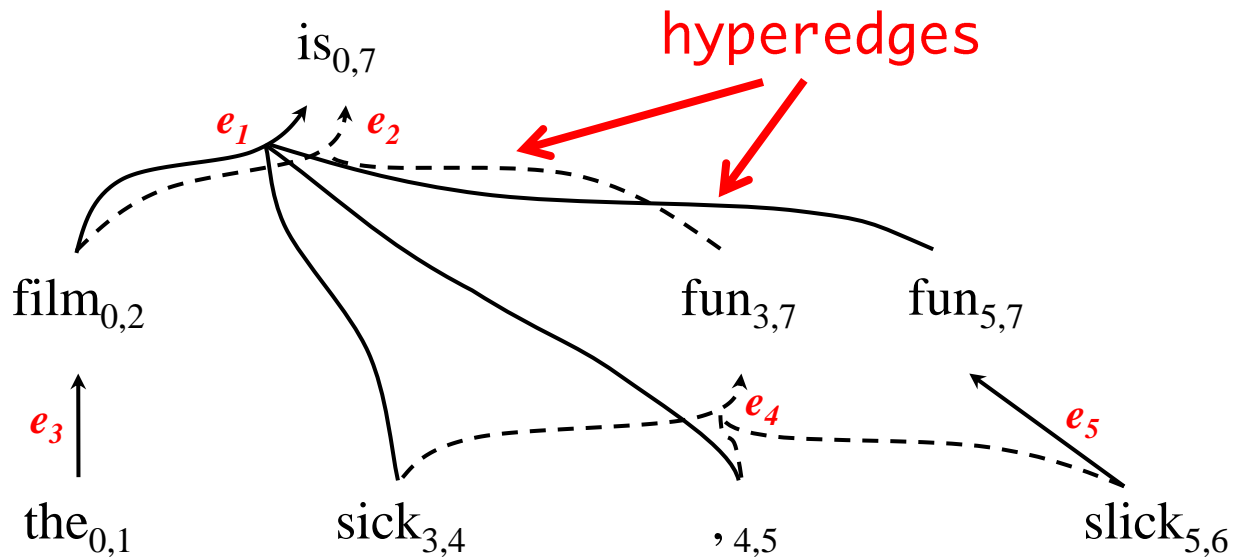
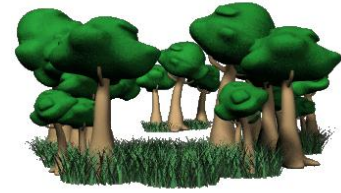




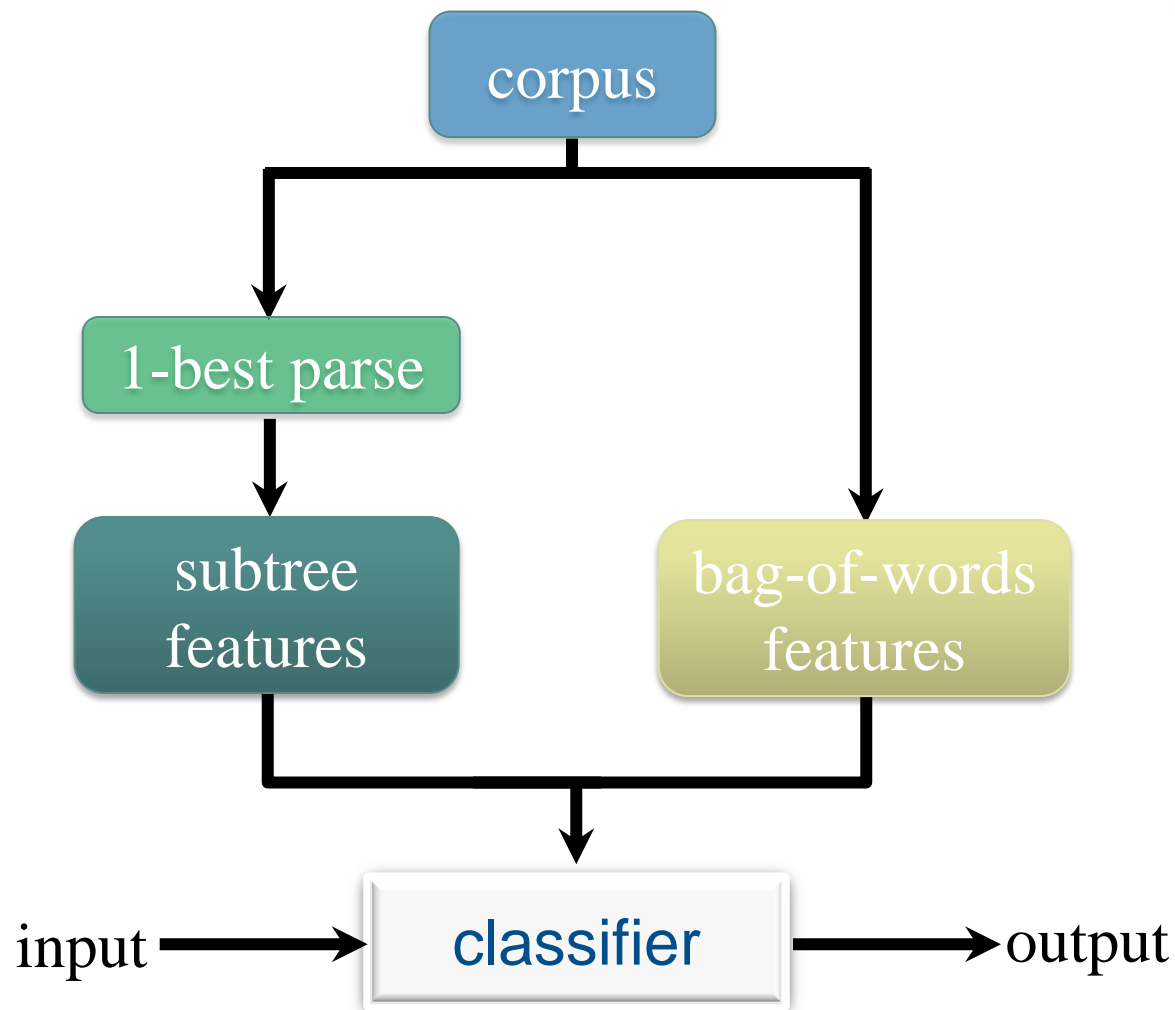
# Dependency Forest



# Dependency Forest

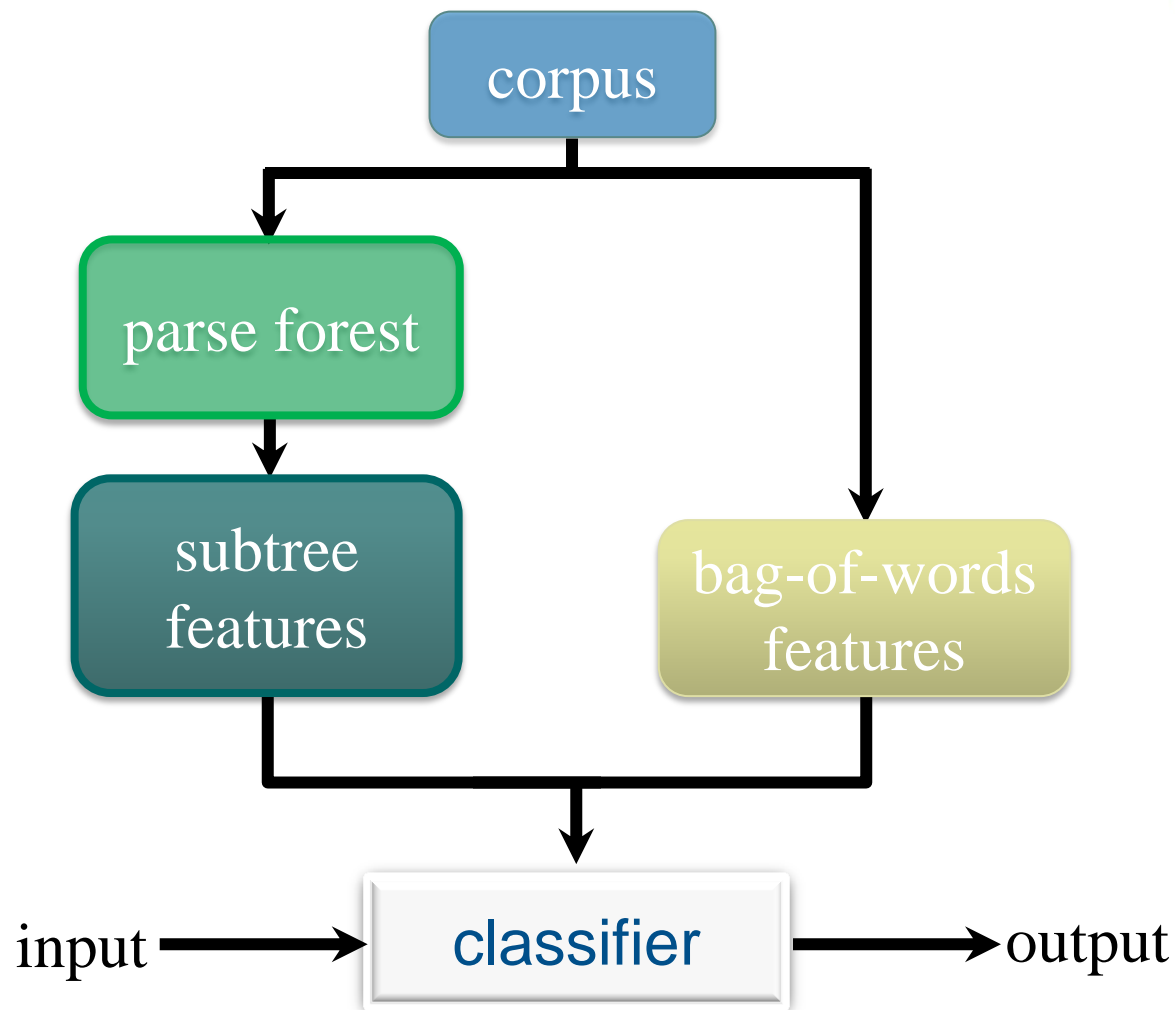


# Pipeline



**tree-based**

# Pipeline

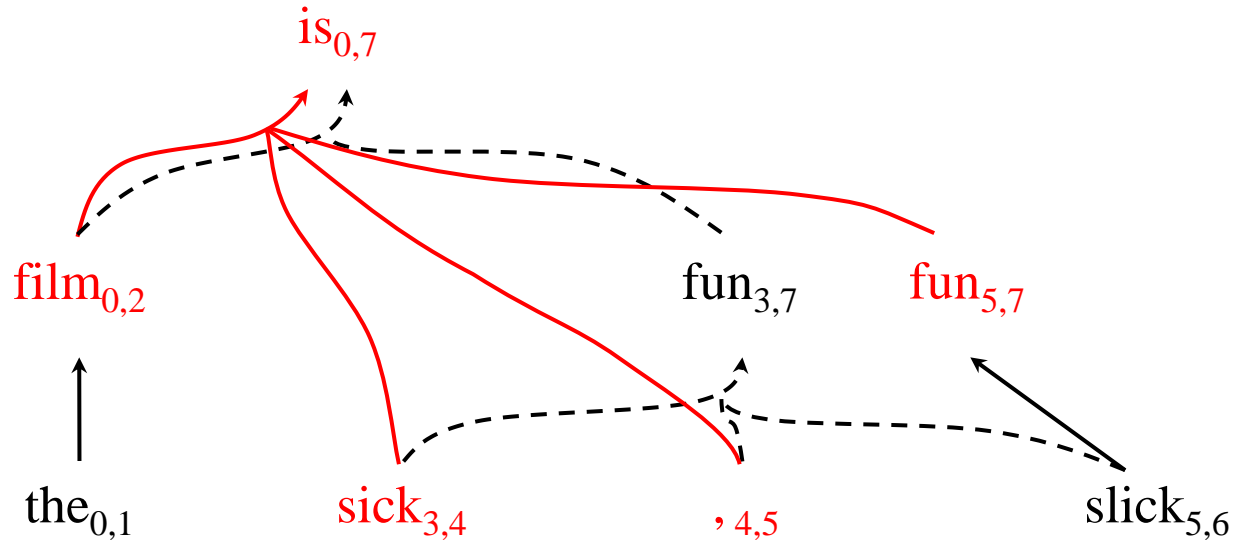


**forest-based**

# Forest-based Feature Extraction



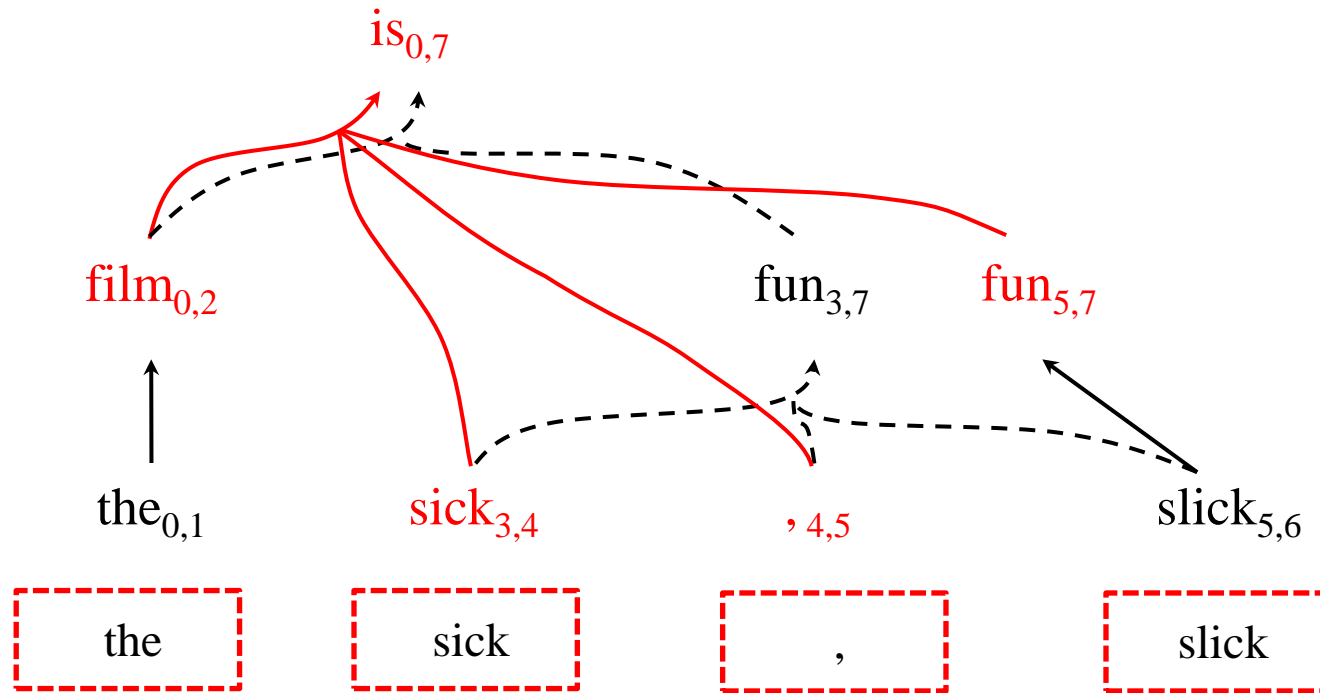
bottom-up style



# Forest-based Feature Extraction



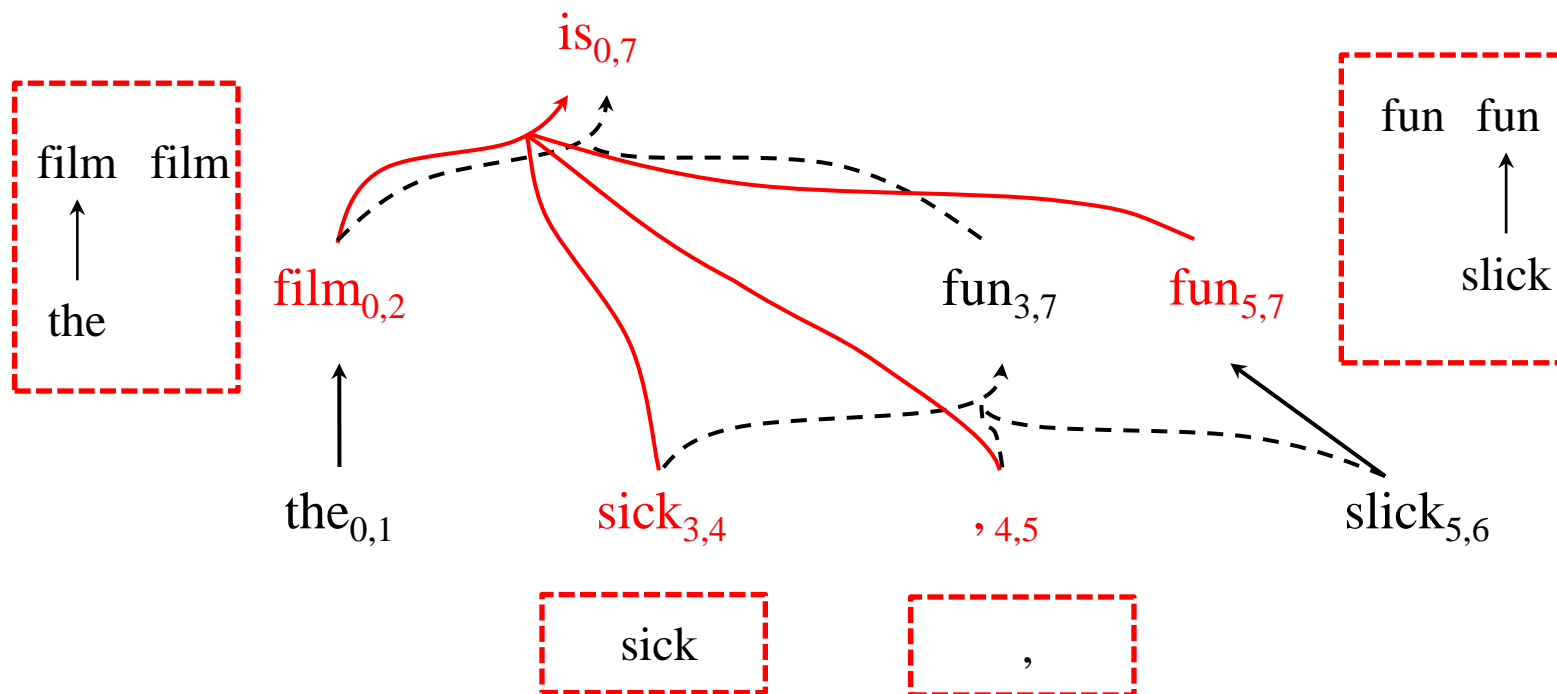
bottom-up style



# Forest-based Feature Extraction



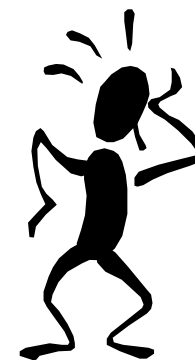
bottom-up style



# Difficulty in Finding Subtrees

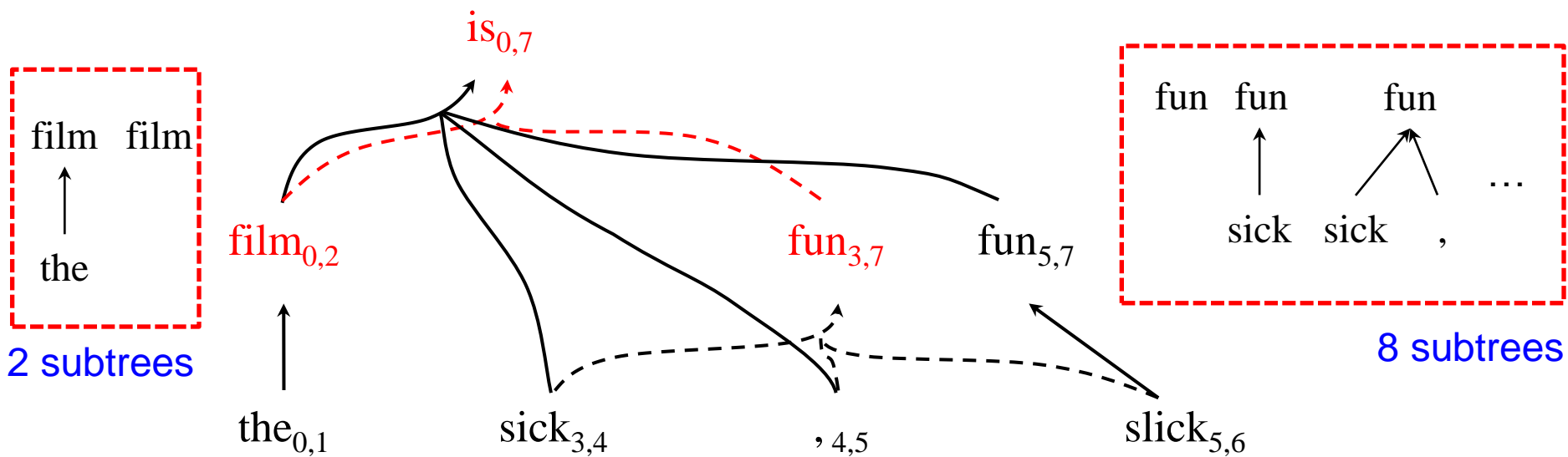


too many choices

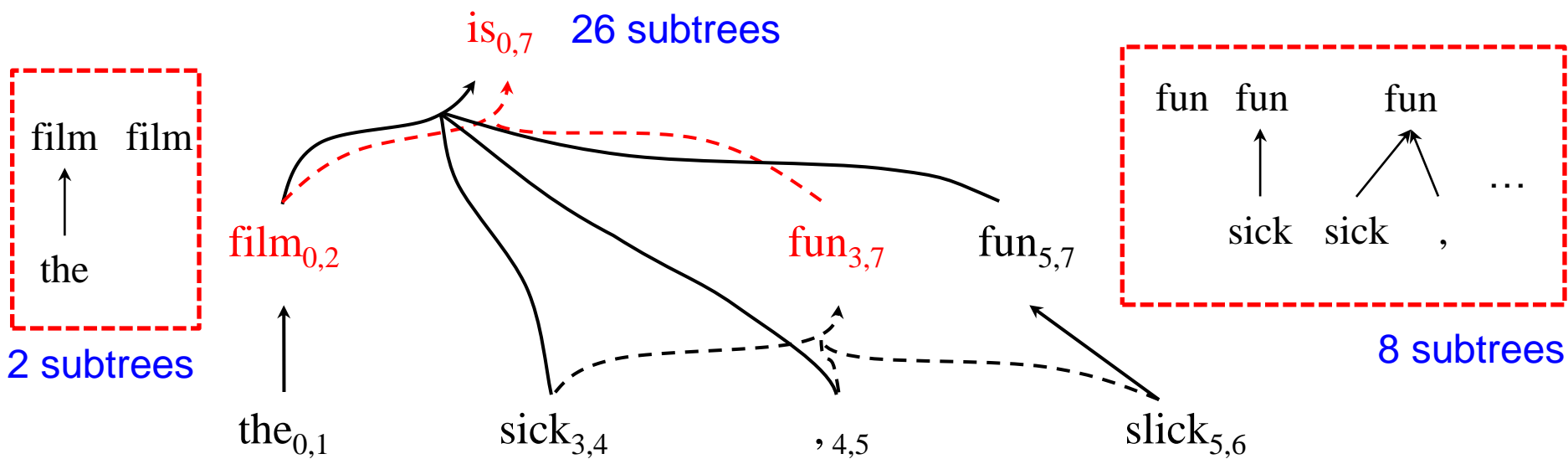




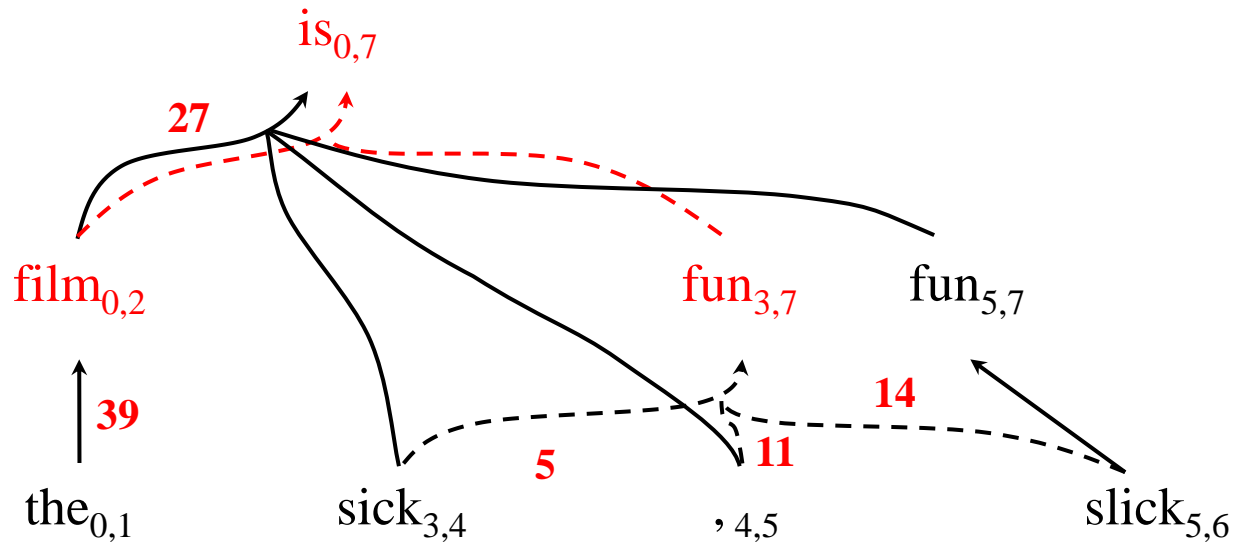
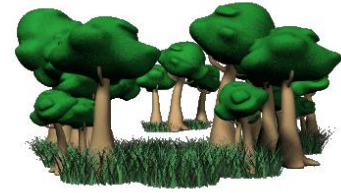
# Forest-based Feature Extraction



# Forest-based Feature Extraction

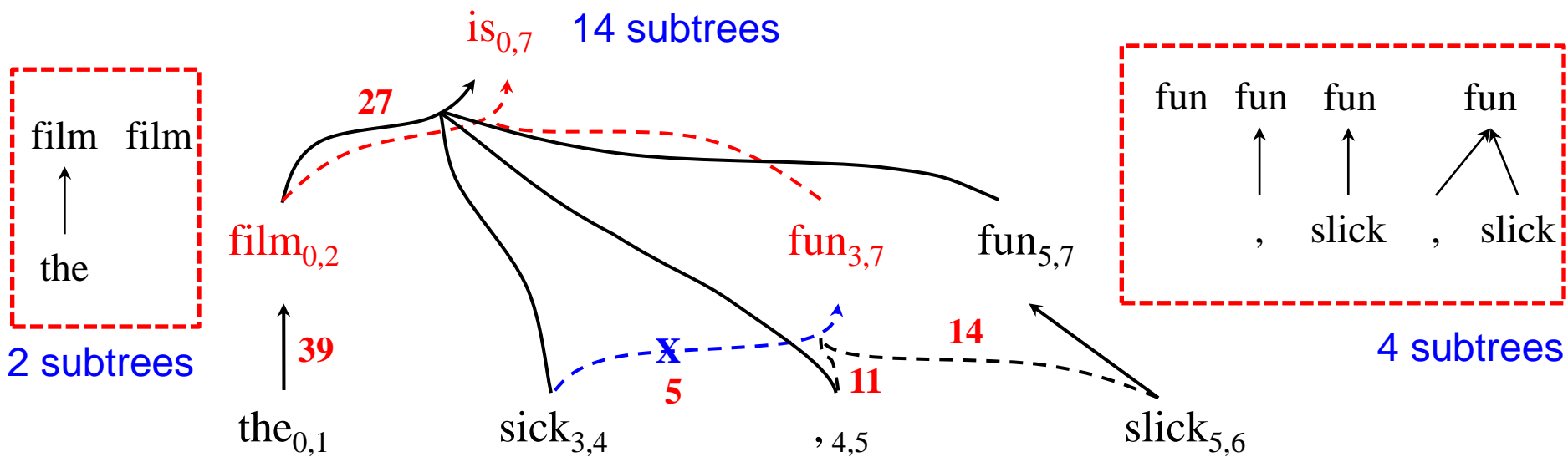


# Forest-based Feature Extraction



gathering the occurrences of edges

# Forest-based Feature Extraction



gathering the occurrences of edges



# Experiments

# Setup



- Movie Review Corpus (Pang et al., 2004)
  - 1000 positive documents and 1000 negative documents
- Unigrams:
  - appear at least 4 times
- Tree-based Subtrees:
  - appear in at least 10 distinct sentences
- Forest-based Subtrees:
  - appear in at least 10 distinct sentences
  - the fractional count should be not lower than 0.01 (Tu et al., 2010)

# Results



Features	Number	Accuracy	Time
<b>Unigram</b>	17,704	86.2	—
<b>Subtree<sub>1–best</sub></b>	12,282	74.2	0.33
<b>Unigram + Subtree<sub>1–best</sub></b>	29,986	90.3 <sup>†</sup>	—
<b>Subtree<sub>100–best</sub></b>	24,006	81.9	35.47
<b>Unigram + Subtree<sub>100–best</sub></b>	41,710	90.2 <sup>†</sup>	—
<b>Subtree<sub>forest</sub></b>	18,968	81.2	6.93
<b>Unigram + Subtree<sub>forest</sub></b>	36,674	<b>91.6<sup>‡</sup></b>	—
Pang et al. [18]	—	87.1	—
Ng et al. [17]	—	90.5	—
Yessenalina et al. [23]	—	91.8	—

# Forest VS Bag-of-words



Features	Number	Accuracy	Time
<b>Unigram</b>	17,704	86.2	—
<b>Subtree<sub>1–best</sub></b>	12,282	74.2	0.33
<b>Unigram + Subtree<sub>1–best</sub></b>	29,986	90.3 <sup>†</sup>	—
<b>Subtree<sub>100–best</sub></b>	24,006	81.9	35.47
<b>Unigram + Subtree<sub>100–best</sub></b>	41,710	90.2 <sup>†</sup>	—
<b>Subtree<sub>forest</sub></b>	18,968	81.2	6.93
<b>Unigram + Subtree<sub>forest</sub></b>	36,674	<b>91.6<sup>‡</sup></b>	—
Pang et al. [18]	—	87.1	—
Ng et al. [17]	—	90.5	—
Yessenalina et al. [23]	—	91.8	—



# Forest VS 1-best tree



Features	Number	Accuracy	Time
<b>Unigram</b>	17,704	86.2	—
<b>Subtree<sub>1-best</sub></b>	12,282	74.2	0.33
<b>Unigram + Subtree<sub>1-best</sub></b>	29,986	90.3 <sup>†</sup>	—
<b>Subtree<sub>100-best</sub></b>	24,006	81.9	35.47
<b>Unigram + Subtree<sub>100-best</sub></b>	41,710	90.2 <sup>†</sup>	—
<b>Subtree<sub>forest</sub></b>	18,968	81.2	6.93
<b>Unigram + Subtree<sub>forest</sub></b>	36,674	<b>91.6<sup>‡</sup></b>	—
Pang et al. [18]	—	87.1	—
Ng et al. [17]	—	90.5	—
Yessenalina et al. [23]	—	91.8	—

# Forest VS $k$ -best trees



Features	Number	Accuracy	Time
<b>Unigram</b>	17,704	86.2	—
<b>Subtree<sub>1-best</sub></b>	12,282	74.2	0.33
<b>Unigram + Subtree<sub>1-best</sub></b>	29,986	90.3 <sup>†</sup>	—
<b>Subtree<sub>100-best</sub></b>	24,006	81.9	35.47
<b>Unigram + Subtree<sub>100-best</sub></b>	41,710	90.2 <sup>†</sup>	—
<b>Subtree<sub>forest</sub></b>	18,968	81.2	6.93
<b>Unigram + Subtree<sub>forest</sub></b>	36,674	<b>91.6<sup>‡</sup></b>	—
Pang et al. [18]	—	87.1	—
Ng et al. [17]	—	90.5	—
Yessenalina et al. [23]	—	91.8	—

# Forest VS State-of-the-art



Features	Number	Accuracy	Time
<b>Unigram</b>	17,704	86.2	—
<b>Subtree<sub>1-best</sub></b>	12,282	74.2	0.33
<b>Unigram + Subtree<sub>1-best</sub></b>	29,986	90.3 <sup>†</sup>	—
<b>Subtree<sub>100-best</sub></b>	24,006	81.9	35.47
<b>Unigram + Subtree<sub>100-best</sub></b>	41,710	90.2 <sup>†</sup>	—
<b>Subtree<sub>forest</sub></b>	18,968	81.2	6.93
<b>Unigram + Subtree<sub>forest</sub></b>	36,674	<b>91.6<sup>‡</sup></b>	—
Pang et al. [18]	—	87.1	—
Ng et al. [17]	—	90.5	—
Yessenalina et al. [23]	—	91.8	—

# Conclusion and Future Work



- dependency forest provides an elegant solution to the problem of *parsing error propagation*
- very simple idea, but works very well in practice
  - 5.4 points in accuracy better than bag-of-words
  - 1.3 points in accuracy better than 1-best trees

Forest offers more alternatives.



**Thank you!**

Thanks to the anonymous reviewers.