# Sentence Simplification with Deep Reinforcement Learning

*Xingxing Zhang and Mirella Lapata*

# What is sentence simplification

Deletion-based sentence compression:
  A man suffered a serious head injury after a morning car crash today .
  A man suffered a injury after a crash .

Sentence simplification:
  1. Delete elements of the original text
  2. Substitute rare words with more common words or phrases
  3. Make syntactically complex structures simpler
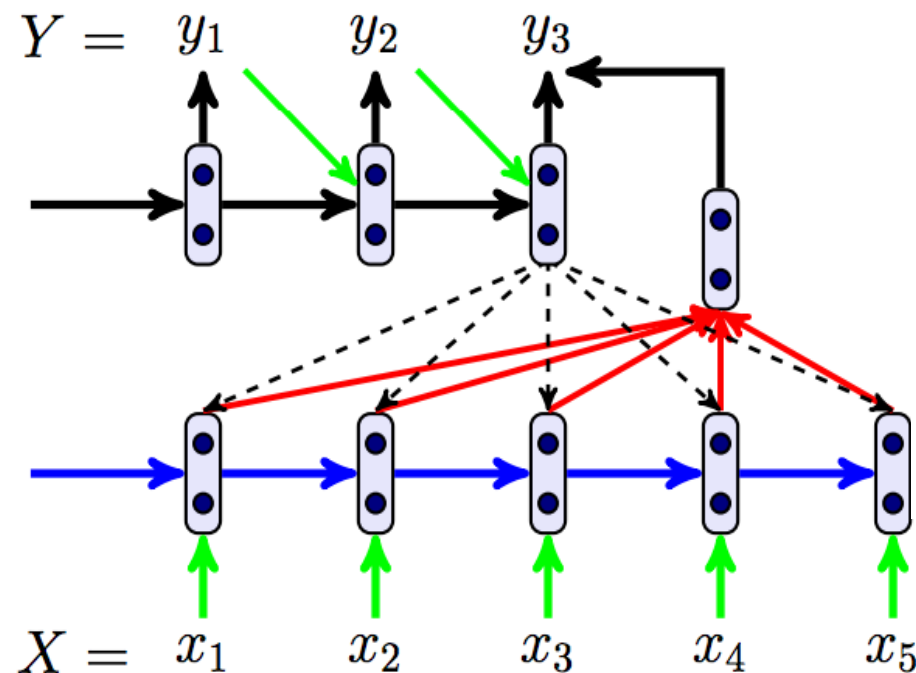
# Application of sentence simplification

- Improve the performance of parsers *(Chandrasekar et al., 1996)*
- Summarizers *(Beigman Klebanov et al., 2004)*
- Semantic role labelers *(Woodsend and Lapata, 2014)*
- Benefit people with low-literacy skills such as children and non-native speakers *(Watanabe et al., 2009)*

# Previous work V.S. Recent work

- Previous works focused on individual aspects of the simplification problem:

  (1) perform syntactic simplification only.
  (2) lexical simplification by substituting difficult words with more common WordNet synonyms or paraphrases.

- Recent works view it as a monolingual text-to-text generation task

# Vanilla Encoder-Decoder with Attention Model

Given a (complex) source sentence: $X = (x_1, x_2, \ldots, x_{|X|})$,
Predict its simplified target: $Y = (y_1, y_2, \ldots, y_{|Y|})$.

$$P(Y|X) = \prod_{t=1}^{|Y|} P(y_t|y_{1:t-1}, X) \qquad (1)$$

⟵ • Minimizing the **negative log-likelihood** of the training source-target pairs

$$P(y_{t+1}|y_{1:t}, X) = \text{softmax}(g(\mathbf{h}_t^T, \mathbf{c}_t)) \qquad (2)$$

where $g(\cdot)$ is a one-hidden-layer neural network with the following parametrization:
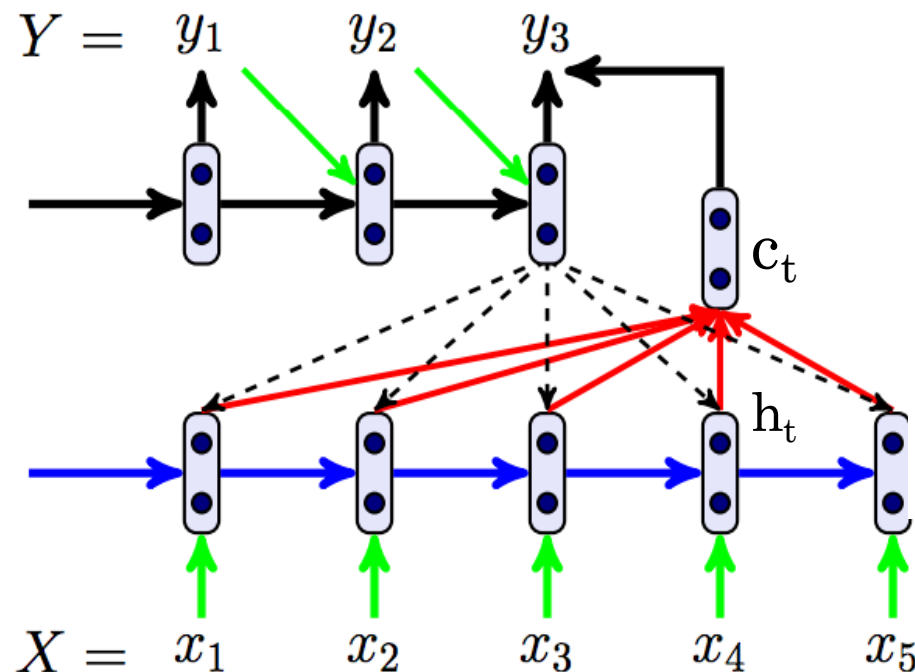
$$g(\mathbf{h}_t^T, \mathbf{c}_t) = \mathbf{W}_o \tanh(\mathbf{U}_h \mathbf{h}_t^T + \mathbf{W}_h \mathbf{c}_t) \qquad (3)$$

where $\mathbf{W}_o \in \mathbb{R}^{|V|\times d}$, $\mathbf{U}_h \in \mathbb{R}^{d\times d}$, and $\mathbf{W}_h \in \mathbb{R}^{d\times d}$; $|V|$ is the output vocabulary size and $d$ the hidden unit size. $\mathbf{h}_t^T$ is the hidden state of the decoder LSTM which summarizes $y_{1:t}$, i.e., what has been generated so far:

$$\mathbf{h}_t^T = \text{LSTM}(y_t, \mathbf{h}_{t-1}^T) \qquad (4)$$

The dynamic context vector $\mathbf{c}_t$ is the weighted sum of the hidden states of the source sentence:
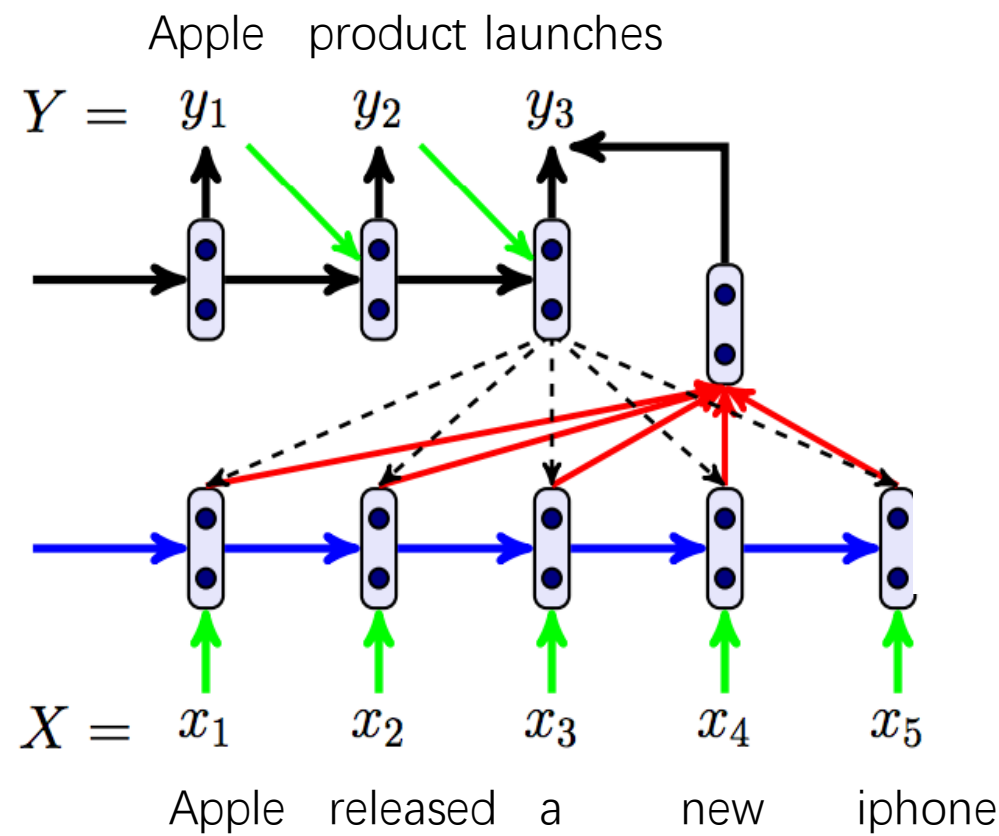
$$\mathbf{c}_t = \sum_{i=1}^{|X|} \alpha_{ti} \mathbf{h}_i^S \qquad (5)$$

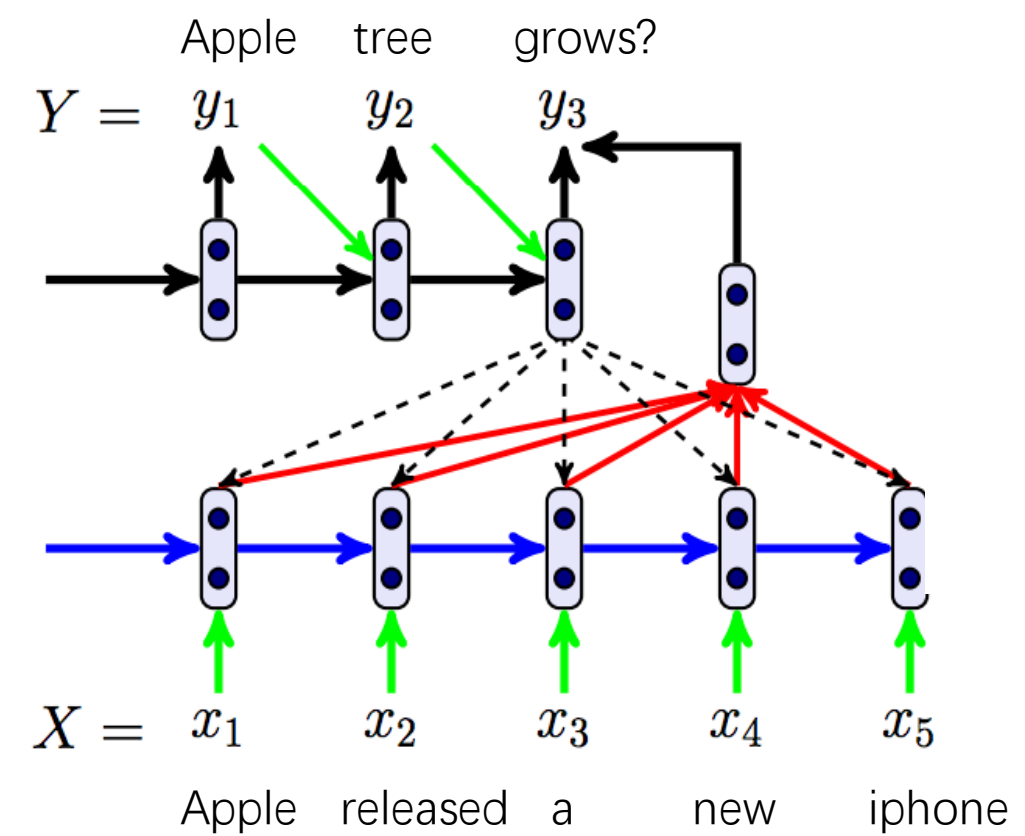# Vanilla seq2seq is not ideal for sentence simplification

- Rewrite operations (e.g., copying, deletion, substitution, word reordering)


- Problem
  - Copy occupied of 73% operation in Newsela dataset, and 83% operation in Wikipedia-based dataset

# Vanilla seq2seq may not be ideal for considering the sentence as a whole



During training

During testing

# Motivation

- To encourage a wider variety of rewrite operations while remaining fluent and faithful to the meaning of the source.

# What leads to a good simplification - Simplicity

- Simplicity: System output Against References and against the Input sentence (SARI) which is the <u>arithmetic average</u> of n-gram precision and recall of three rewrite operations: <u>addition</u>, <u>copying</u>, and <u>deletion</u>.
  - reward the addition operations **where system output was not in the input but occurred in the references**.

- What the paper used is: SARI score(e.x. 0.7594)

# What leads to a good simplification - Relevance

- **Relevance**: while encouraging changes, generated sentence should preserve the meaning of the source.

  - Cosine similarity between original one and simply one.

$$r^R = \cos(\mathbf{q}_X, \mathbf{q}_{\hat{Y}}) = \frac{\mathbf{q}_X \cdot \mathbf{q}_{\hat{Y}}}{\|\mathbf{q}_X\| \, \|\mathbf{q}_{\hat{Y}}\|}$$

  - $q_x$ and $q_y$ are vector representation of source and target.

# What leads to a good simplification - Fluency

- Fluency: be readable & be grammatical

  - LSTM language model trained on simple sentences

$$r^F = \exp \left( \underbrace{\frac{1}{|\hat{Y}|} \sum_{i=1}^{|\hat{Y}|} \log P_{LM}(\hat{y}_i | \hat{y}_{0:i-1})}_{\text{perplexity}} \right)$$

We take the exponential of Y's perplexity to ensure that $r^F \in [0,1]$

# Put them together

simplicity, relevance, and fluency:

$$r(\hat{Y}) = \lambda^S r^S + \lambda^R r^R + \lambda^F r^F$$

where $\lambda^S, \lambda^R, \lambda^F \in [0, 1]$; $r(\hat{Y})$ is a shorthand for $r(X, Y, \hat{Y})$ where $X$ is the source, $Y$ the reference (or target), and $\hat{Y}$ the system output.

The reward r(Yˆ) for system output Yˆ is the weighted sum of the three components
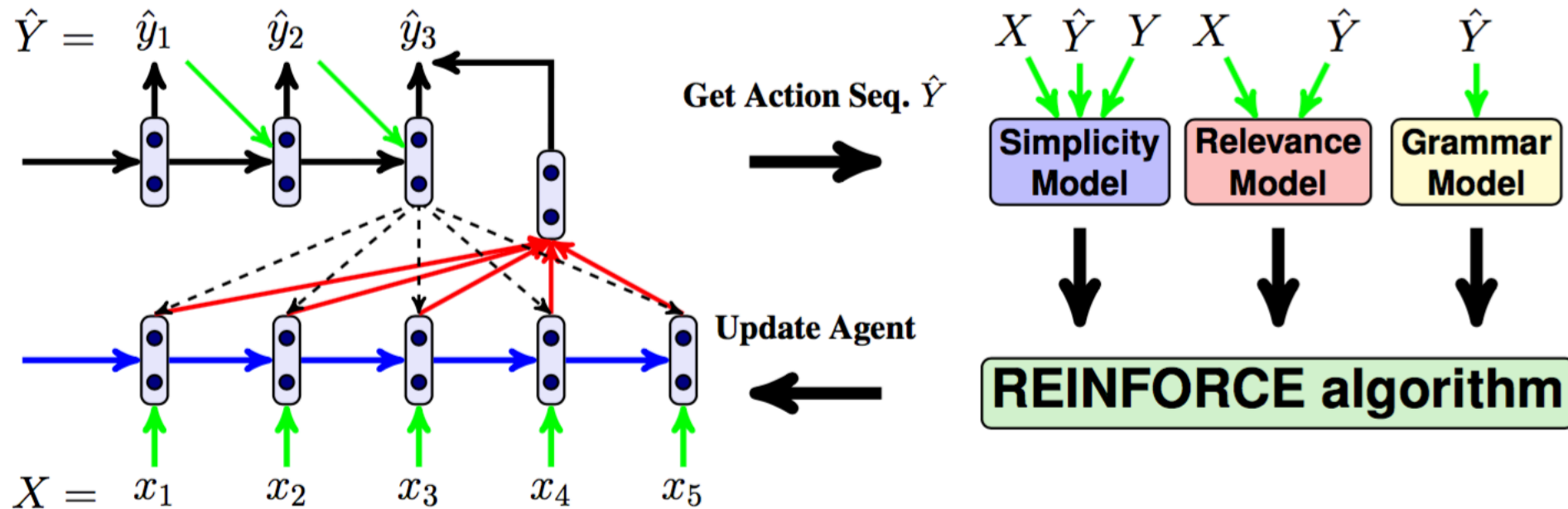
# Put them together



Figure 1: Deep reinforcement learning simplification model. $X$ is the complex sentence, $Y$ the reference (simple) sentence and $\hat{Y}$ the action sequence (simplification) produced by the encoder-decoder model.

# Datasets

Automatic
alignment

- Parallel corpus$_1$ *WikiSmall*
  - It contains automatically aligned complex and simple sentences from the ordinary and simple English Wikipedias.
  - train/val/test  89,042/205/100

- Parallel corpus$_2$ *WikiLarge*
  - It contains 8 (reference) simplifications for 2,359 sentences partitioned into 2,000 for development and 359 for testing.
  - train/val/test  296,402/2,000/359

- Parallel corpus$_3$ Newsela
  - It consists of 1,130 news articles, each rewritten four times by professional editors for children at different grade levels
  - train/val/test  94,208/1,129/1,076

professional
editors

# Automatic Evaluation

1. Bilingual Evaluation Understudy BLEU[blɛː]
   - To assess the degree to which generated simplifications differed from gold standard references

2. Flesch-Kincaid Grade Level (FKGL) score
   - To measure the simplicity of the output (lower FKGL implies simpler output)

3. System output Against References and against the Input sentence (SARI)
   - To evaluate the quality of the output by comparing it against the source and reference simplifications

# Human Evaluation

Native English speakers were asked to rate simplifications on three dimensions:

- *Fluency* (is the output grammatical and well formed?),
- *Adequacy* (to what extent is the meaning expressed in the original sentence preserved in the output?)
- *Simplicity* (is the output simpler than the original sentence?)

All ratings were obtained using a five point Likert scale.

# Results

| Newsela | BLEU | FKGL | SARI |
|---------|------|------|------|
| PBMT-R | 18.19 | 7.59 | 15.77 |
| Hybrid | 14.46 | **4.01** | **30.00** |
| EncDecA | 21.70 | 5.11 | 24.12 |
| DRESS | 23.21 | 4.13 | 27.37 |
| DRESS-LS | **24.30** | 4.21 | 26.63 |

| WikiSmall | BLEU | FKGL | SARI |
|-----------|------|------|------|
| PBMT-R | 46.31 | 11.42 | 15.97 |
| Hybrid | **53.94** | 9.20 | **30.46** |
| EncDecA | 47.93 | 11.35 | 13.61 |
| DRESS | 34.53 | **7.48** | 27.48 |
| DRESS-LS | 36.32 | 7.55 | 27.24 |

| WikiLarge | BLEU | FKGL | SARI |
|-----------|------|------|------|
| PBMT-R | 81.11 | 8.33 | 38.56 |
| Hybrid | 48.97 | **4.56** | 31.40 |
| SBMT-SARI | 73.08 | 7.29 | **39.96** |
| EncDecA | **88.85** | 8.41 | 35.66 |
| DRESS | 77.18 | 6.58 | 37.08 |
| DRESS-LS | 80.12 | 6.62 | 37.27 |

Table 1: Automatic evaluation on Newsela, WikiSmall, and WikiLarge test sets.

| Newsela | Fluency | Adequacy | Simplicity | All |
|---------|---------|----------|------------|-----|
| PBMT-R | 3.56 | **3.58**** | 2.09** | 3.08** |
| Hybrid | 2.70** | 2.51** | 2.99 | 2.73** |
| EncDecA | 3.63 | 2.99 | 2.56** | 3.06** |
| DRESS | 3.65 | 2.94 | **3.10** | 3.23 |
| DRESS-LS | **3.71** | 3.07 | 3.04 | **3.28** |
| Reference | 3.90 | 2.81** | 3.42** | 3.38 |

| WikiSmall | Fluency | Adequacy | Simplicity | All |
|-----------|---------|----------|------------|-----|
| PBMT-R | 3.91 | **3.74**** | 2.80** | 3.48* |
| Hybrid | 3.26** | 3.42 | 2.82** | 3.17** |
| DRESS-LS | **3.92** | 3.36 | **3.55** | **3.61** |
| Reference | 3.74* | 3.34 | 3.13** | 3.41** |

| WikiLarge | Fluency | Adequacy | Simplicity | All |
|-----------|---------|----------|------------|-----|
| PBMT-R | 3.68 | **3.63*** | 2.70** | 3.34* |
| Hybrid | 2.60** | 2.42** | **3.52** | 2.85** |
| SBMT-SARI | 3.34** | 3.51* | 2.77** | 3.21** |
| DRESS-LS | **3.70** | 3.28 | 3.42 | **3.46** |
| Reference | 3.79 | 3.72** | 2.86** | 3.46 |

Table 2: Mean ratings elicited by humans on Newsela, WikiSmall, and WkiLarge test sets. Ratings significantly different from DRESS-LS are marked with * ($p < 0.05$) and ** ($p < 0.01$). Significance tests were performed using a student $t$-test.

# System output for example sentence

| | |
|---|---|
| Complex | There's just one major hitch: the primary purpose of education is to develop citizens with a wide variety of skills. |
| Reference | The purpose of education is to develop a wide range of skills. |
| PBMT-R | It's just one major hitch: the purpose of education is to **make people** with a wide variety of skills. |
| Hybrid | one hitch the purpose is to develop citizens. |
| EncDecA | The **key** of education is to develop **people** with a wide variety of skills. |
| DRESS | There's just one major hitch: the **main goal** of education is to develop **people** with **lots of** skills. |
| DRESS-LS | There's just one major hitch: the **main goal** of education is to develop citizens with **lots of** skills. |

*Substitutions are shown in bold.*

# Conclusion

- Definition of the reward function is the key.

- Take the evaluation metric itself as optimization objective (like SARI score)
  - evaluation is thus important

- …