# CopyNET - Incorporating Copying Mechanism in Seq2Seq Learning

Jiatao Gu, Zhengdong Lu, Hang Li, Victor O.K. Li 2016

Louis (Yiqing) Luo

Aug 3rd, 2018

# Motivation

- Problem: Copying in Seq2Seq
  - certain segments in the input sequence are selectively replicated in the output sequence
  - eg. humans tend to repeat ntity names or even long phrases in conversation

---

I: Hello Jack, my name is Chandralekha.

R: Nice to meet you, Chandralekha.

---

I: This new guy doesn't perform exactly as we expected.

R: What do you mean by "doesn't perform exactly as we expected"?
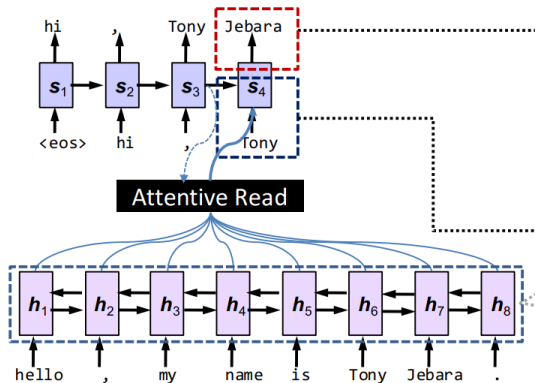
---

# Motivation

- Problem: Copying in Seq2Seq
    - certain segments in the input sequence are selectively replicated in the output sequence
    - eg. humans tend to repeat entity names or even long phrases in conversation

Proposed Solution:

CopyNET = RNN Encoder & Decoder + Copying Mechanism

# RNNSearch: RNN Encoder and Decoder

- typically used in Seq2Seq learning



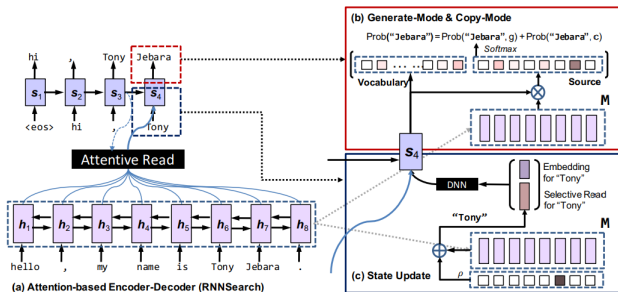**(a) Attention-based Encoder-Decoder (RNNSearch)**

# CopyNET

Encoder same. Decoder Differences:

1. **Prediction:** COPYNET predicts words based on a mixed probabilistic model of two modes
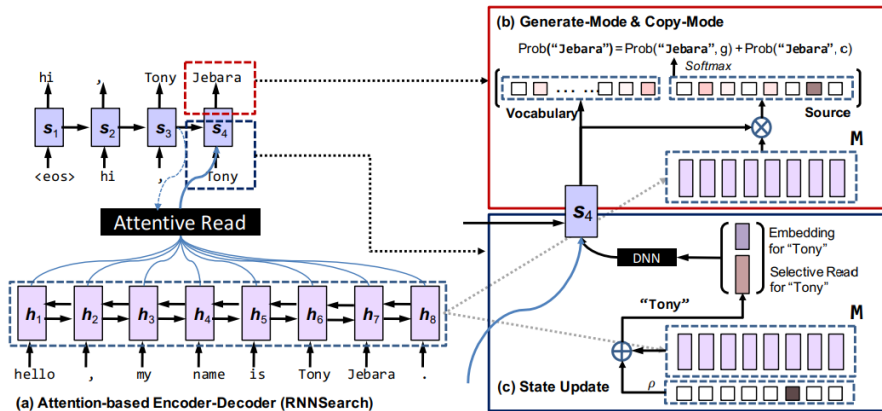2. **State Update:** uses not only its word-embedding but also its corresponding location-specific hidden state in M

Encoder same. Decoder Differences:

1. **Prediction:** COPYNET predicts words based on a mixed probabilistic model of two modes
2. **State Update:** uses not only its word-embedding but also its corresponding location-specific hidden state in M



(a) Attention-based Encoder-Decoder (RNNSearch)

# Model



**(b) Generate-Mode & Copy-Mode**

$\text{Prob}(\text{"Jebara"}) = \text{Prob}(\text{"Jebara"}, g) + \text{Prob}(\text{"Jebara"}, c)$

*Softmax*

Vocabulary | Source

$M$

**(a) Attention-based Encoder-Decoder (RNNSearch)**

Attentive Read

Embedding for "Tony"

Selective Read for "Tony"

DNN

"Tony"

$M$

**(c) State Update**

- For vocabulary $V$ and unique words in source sequence $X$
- Instance-specific Vocabulary for source X is $V \cup UNK \cup X$

Given the decoder RNN state $\mathbf{s}_t$ at time t together with $\mathbf{M}$, the probability of generating any target word $y_t$, is given by the mixture of probabilities as follows:

$$p(y_t|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M}) = p(y_t, \mathbf{g}|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M})$$
$$+ p(y_t, \mathbf{c}|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M}) \quad (4)$$

# Equations Unrolled:

$$p(y_t|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M}) = p(y_t, \mathbf{g}|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M})$$
$$+ p(y_t, \mathbf{c}|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M}) \quad (4)$$

g: generative mode
$$p(y_t, \mathbf{g}|\cdot) = \begin{cases} \frac{1}{Z} e^{\psi_g(y_t)}, & y_t \in \mathcal{V} \\ 0, & y_t \in \mathcal{X} \cap \bar{\mathcal{V}} \quad (5) \\ \frac{1}{Z} e^{\psi_g(\text{UNK})} & y_t \notin \mathcal{V} \cup \mathcal{X} \end{cases}$$

c: copy mode
$$p(y_t, \mathbf{c}|\cdot) = \begin{cases} \frac{1}{Z} \sum_{j:x_j=y_t} e^{\psi_c(x_j)}, & y_t \in \mathcal{X} \quad (6) \\ 0 & \text{otherwise} \end{cases}$$

$$\psi_g(y_t = v_i) = \mathbf{v}_i^\top \mathbf{W}_o \mathbf{s}_t, \quad v_i \in \mathcal{V} \cup \text{UNK}$$
$$\psi_c(y_t = x_j) = \sigma\left(\mathbf{h}_j^\top \mathbf{W}_c\right) \mathbf{s}_t, \quad x_j \in \mathcal{X}$$

Normalizing constant $Z = \sum_{v \in \mathcal{V} \cup \{\text{UNK}\}} e^{\psi_g(v)} + \sum_{x \in X} e^{\psi_c(x)}$.



$\frac{1}{Z} \sum_{x_j} \exp[\psi_c(x_j)] \mid x_j = y_t$

$\frac{1}{Z} \exp[\psi_g(v_i)] \mid v_i = y_t$

$\frac{1}{Z}\left(\sum_{x_j} [\psi_c(x_j)] + \exp[\psi_g(v_i)]\right) \mid x_j = y_t, v_i = y_t$

$\frac{1}{Z} \exp[\psi_g(\text{unk})]$
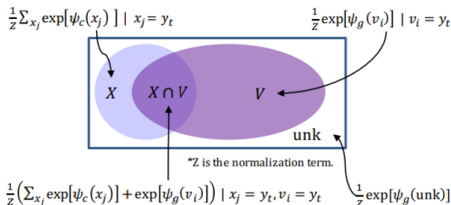
*Z is the normalization term.

Figure 2: The illustration of the decoding probability $p(y_t|\cdot)$ as a 4-class classifier.

# Decoder: 2. State Update

- normally $\mathbf{s}_t$ is updated by $\mathbf{s}_{t-1}, y_{t-1}, and \mathbf{c}_t$
- with CopyNET: the $y_{t-1}$ in $y_{t-1} \to \mathbf{s}_t$ is replaced with:

$$[\mathbf{e}(y_{t-1}); \zeta(y_{t-1})]^\top$$

where $\quad \zeta(y_{t-1}) = \sum_{\tau=1}^{T_S} \rho_{t\tau} \mathbf{h}_\tau$

$$\rho_{t\tau} = \begin{cases} \frac{1}{K} p(x_\tau, \mathbf{c} | \mathbf{s}_{t-1}, \mathbf{M}), & x_\tau = y_{t-1} \\ 0 & \text{otherwise} \end{cases}$$

$\mathsf{K} = \sum_{\tau' : x_{\tau'} = y_{t-1}} p(x_{\tau'}, c | \mathbf{s}_{t-1}, \mathbf{M})$

where $\mathbf{e}(y_t 1)$ is the word embedding associated with $y_t 1$, while $\zeta(y_t 1)$ is the weighted sum of hidden states in $\mathbf{M}$ corresponding to $y_t$

# Loss function and Updating

$$\mathcal{L} = -\frac{1}{N} \sum_{k=1}^{N} \sum_{t=1}^{T} \log \left[ p(y_t^{(k)} | y_{<t}^{(k)}, X^{(k)}) \right]$$

where source sequence $= X^{(N)}$ and target sequence $= Y^{(N)}$

- The network can learn to coordinate the two modes from data
  - if a target word in the source sequence, the copy-mode will contribute to the mixture model, and the gradient will more or less encourage the copy-mode; otherwise, the copy-mode is discouraged due to the competition from the shared normalization term Z

# Experiment

1. A synthetic dataset on with simple patterns;
2. A real-world task on text summarization;
3. A dataset for simple single-turn dialogues.

# Experiments - Synthetic Dataset

Each rule can further produce a number of instances by replacing the variables with randomly generated subsequences (1 to 15 symbols) from the same vocabulary

| Rule-type | Examples (e.g. $\mathbf{x} =$ i h k, $\mathbf{y} =$ j c) |
|---|---|
| $\mathbf{x} \rightarrow \emptyset$ | a b c d $\mathbf{x}$ e f $\rightarrow$ c d g |
| $\mathbf{x} \rightarrow \mathbf{x}$ | a b c d $\mathbf{x}$ e f $\rightarrow$ c d $\mathbf{x}$ g |
| $\mathbf{x} \rightarrow \mathbf{x}\,\mathbf{x}$ | a b c d $\mathbf{x}$ e f $\rightarrow$ $\mathbf{x}$ d $\mathbf{x}$ g |
| $\mathbf{x}\,\mathbf{y} \rightarrow \mathbf{x}$ | a b $\mathbf{y}$ d $\mathbf{x}$ e f $\rightarrow$ $\mathbf{x}$ d i g |
| $\mathbf{x}\,\mathbf{y} \rightarrow \mathbf{x}\,\mathbf{y}$ | a b $\mathbf{y}$ d $\mathbf{x}$ e f $\rightarrow$ $\mathbf{x}$ d $\mathbf{y}$ g |

| Rule-type | $x \to \emptyset$ | $x \to x$ | $x \to xx$ | $xy \to x$ | $xy \to xy$ |
|---|---|---|---|---|---|
| Enc-Dec | **100** | 3.3 | 1.5 | 2.9 | 0.0 |
| RNNSearch | 99.0 | 69.4 | 22.3 | 40.7 | 2.6 |
| COPYNET | 97.3 | **93.7** | **98.3** | **68.2** | **77.5** |

Table 1: The test accuracy (%) on synthetic data.

- Encoder-Decoder (no Attention) $\to$ difficulty of representing a long sequence with very high fidelity
- RNNSearch (with Attention) $\to$ attention alone seems inadequate for handling the case where strict replication is needed

Automatic text summarization aims to find a condensed representation which can capture the core meaning of the original document

- Dataset: LCSTS dataset (Hu et al., 2015), a large scale dataset for short text summarization in form of (short news, summary).
- model tried on character ($+C$) and word ($+W$)
- ROUGE-N: Overlap of N-grams between the system and reference summaries.
- ROUGE - LCN: measures longest matching sequence of words using longest common subsequence.

| Models | | ROUGE scores on LCSTS (%) | | |
|---|---|---|---|---|
| | | R-1 | R-2 | R-L |
| RNN | +C | 21.5 | 8.9 | 18.6 |
| (Hu et al., 2015) | +W | 17.7 | 8.5 | 15.8 |
| RNN context | +C | 29.9 | 17.4 | 27.2 |
| (Hu et al., 2015) | +W | 26.8 | 16.1 | 24.1 |
| COPYNET | +C | **34.4** | **21.6** | **31.3** |
| | +W | **35.0** | **22.3** | **32.0** |

Table 3: Testing performance of LCSTS, where

# Experiments - Text Summarization



Figure 4: Examples of CopyNet on LCSTS compared with RNN context. Word segmentation is applied on the input, where OOV words are underlined. The highlighted words (with different colors) are those words with copy-mode probability higher than the generate-mode. We also provide literal

# Experiments - Text Summarization

1. most words are from copy-mode, but the summary is usually still fluent;

2. COPYNET tends to cover consecutive words in the original document, but it often puts together segments far away from each other, indicating a sophisticated coordination of content-based addressing and location-based addressing;

3. COPYNET handles OOV words really well: it can generate acceptable summary for document with many OOVs, and even the summary itself often contains many OOV words

# Experiments - Single-turn Dialogue

- Dataset built via a simple dialogue dataset based on the following three instructions
  1. Dialogue instances are collected from Baidu Tieba3 with some coverage of conversations of real life e.g., greeting and sports, etc.
  2. patterns with slots like

     hi, my name is x   hi, x

     are mined from the set, with possibly multiple responding patterns to one input.
  3. Similar with the synthetic dataset, we enlarge the dataset by filling the slots with suitable subsequence (e.g. name entities, dates, etc.)
     - Created 2 datasets: DS-I and DS-II
     - the filled substrings for training and testing in DS-II have no overlaps, while in DS-I they are sampled from the same pool

| Models | DS-I (%) | | DS-II (%) | |
|---|---|---|---|---|
| | Top1 | Top10 | Top1 | Top10 |
| RNNSearch | 44.1 | 57.7 | 13.5 | 15.9 |
| COPYNET | **61.2** | **71.0** | **50.5** | **64.8** |

- Both models estimate respectively the chance of the top-1 or one of top-10 (from beam search) matching the golden.

Figure 5: Examples from the testing set of DS-II shown as the input text and golden, with the outputs of RNNSearch and CopyNet. Words in red rectangles are unseen in the training set. The highlighted words (with different colors) are those words with copy-mode probability higher than the generate-mode. Green cirles (meaning correct) and red cross (meaning incorrect) are given based on human judgment on whether the response is appropriate.