# Presentation 2:
# Long Document Summarization

CS4624 Multimedia/Hypertext/Information Access

Team: Junjie Cheng

Instructor: Edward A. Fox

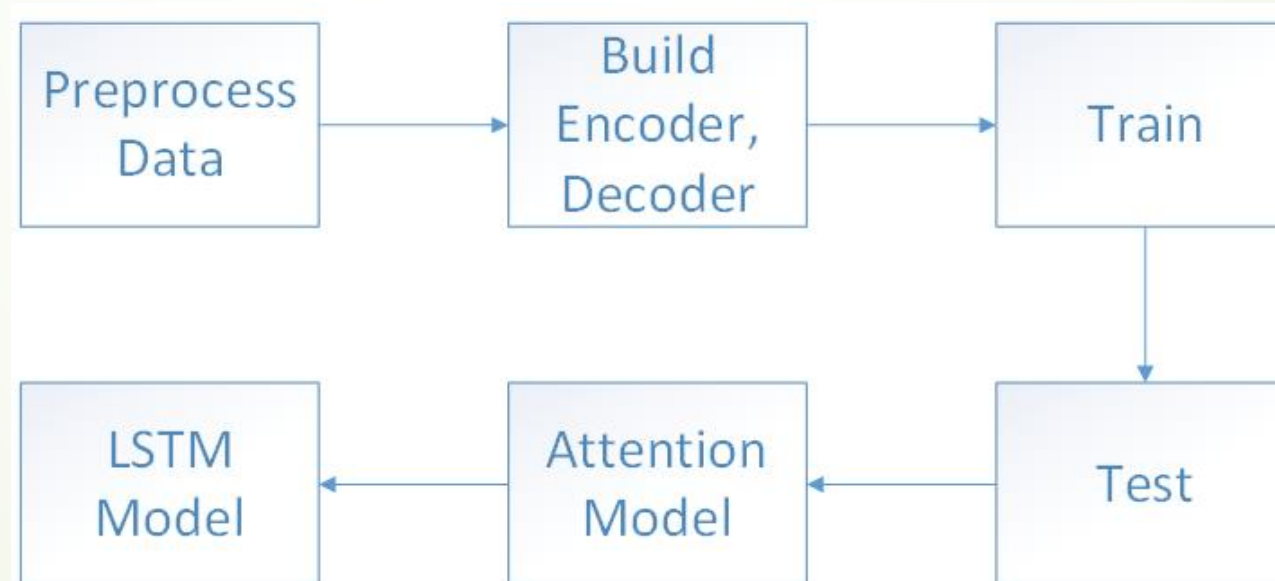Virginia Tech, Blacksburg VA 24061, Mar 13, 2018

# Outline

- Project Design
- Current Status
  - Data Preprocess
  - Encoder/Decoder
  - Training
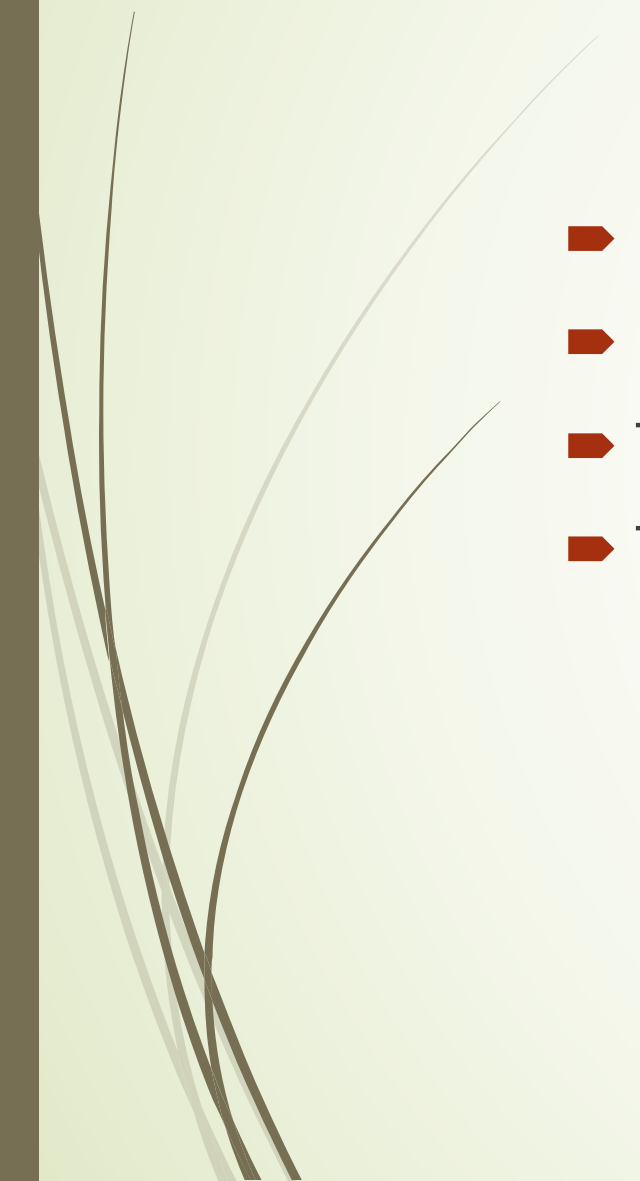  - Testing
- Future Plan
- Acknowledgements and References

# Project Design

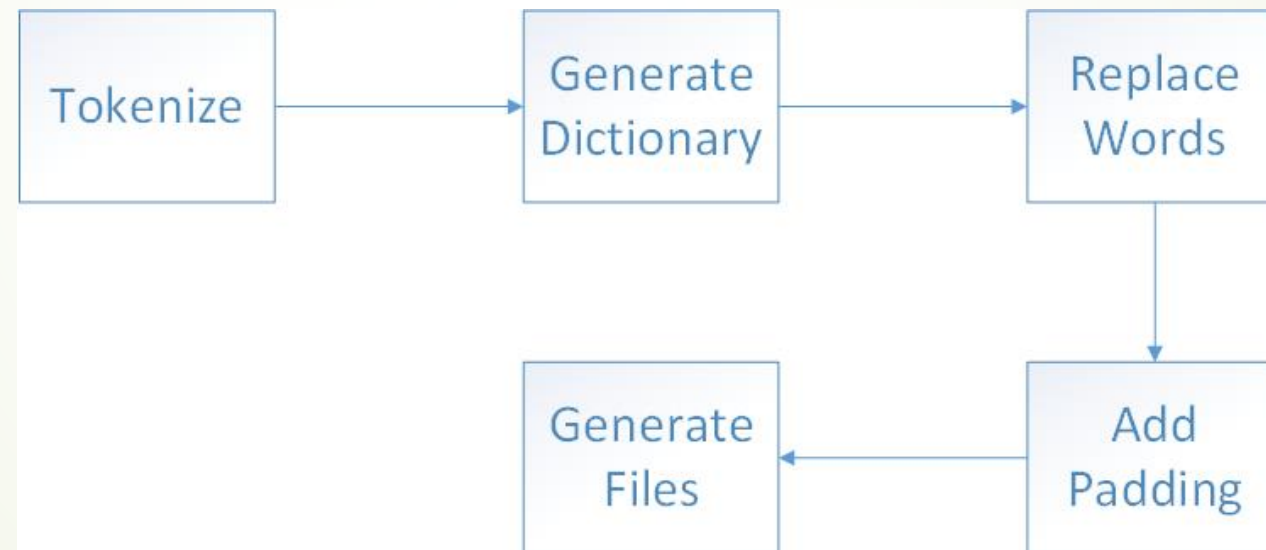■ **Purpose**: Generate abstract from long document by deep learning.

# Current Status

- Data Preprocess

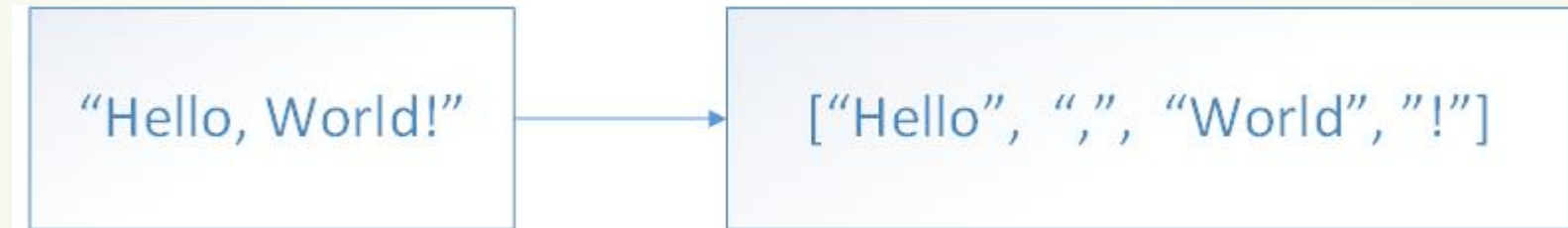- Encoder/Decoder

- Training

- Testing

# Data Preprocess

- 300,000 articles from CNN/Dailymail

# Data Preprocess: Tokenize

- Tokenize words, symbols and numbers

"Hello, World!" → ["Hello", ",", "World", "!"]

# Data Preprocess: Generate Dictionary

- Generate a dictionary includes all words and symbols
- For example:
  - "Hello" → 4
  - "," → 5
- Also include:
  - "<Padding>" → 0
  - "<Start of Sentence>" → 1
  - "<End of Sentence>" → 2
  - "<Unknown Word>" → 3

# Data Preprocess:
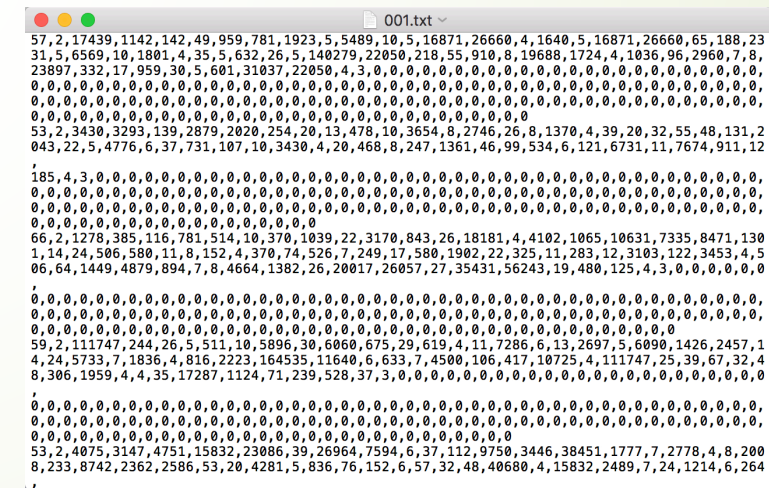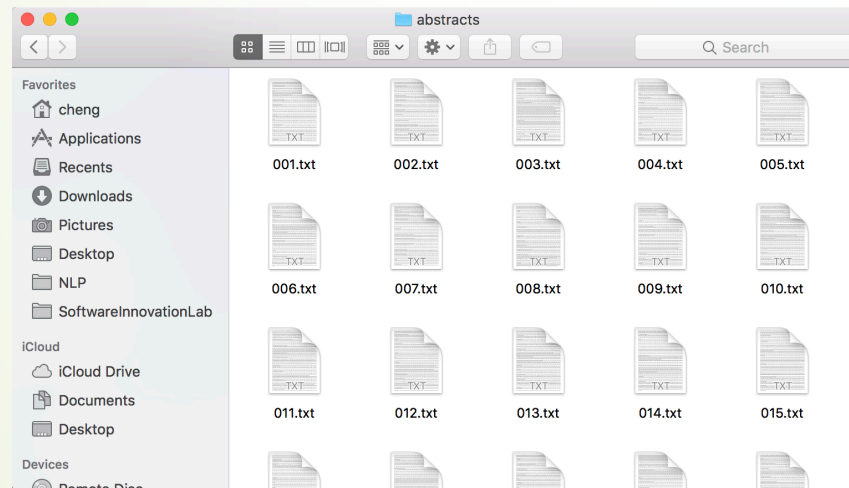# Replace Words and Add Padding

- Replace words by ID
- Add <SOS>, <EOS>, and <Padding>
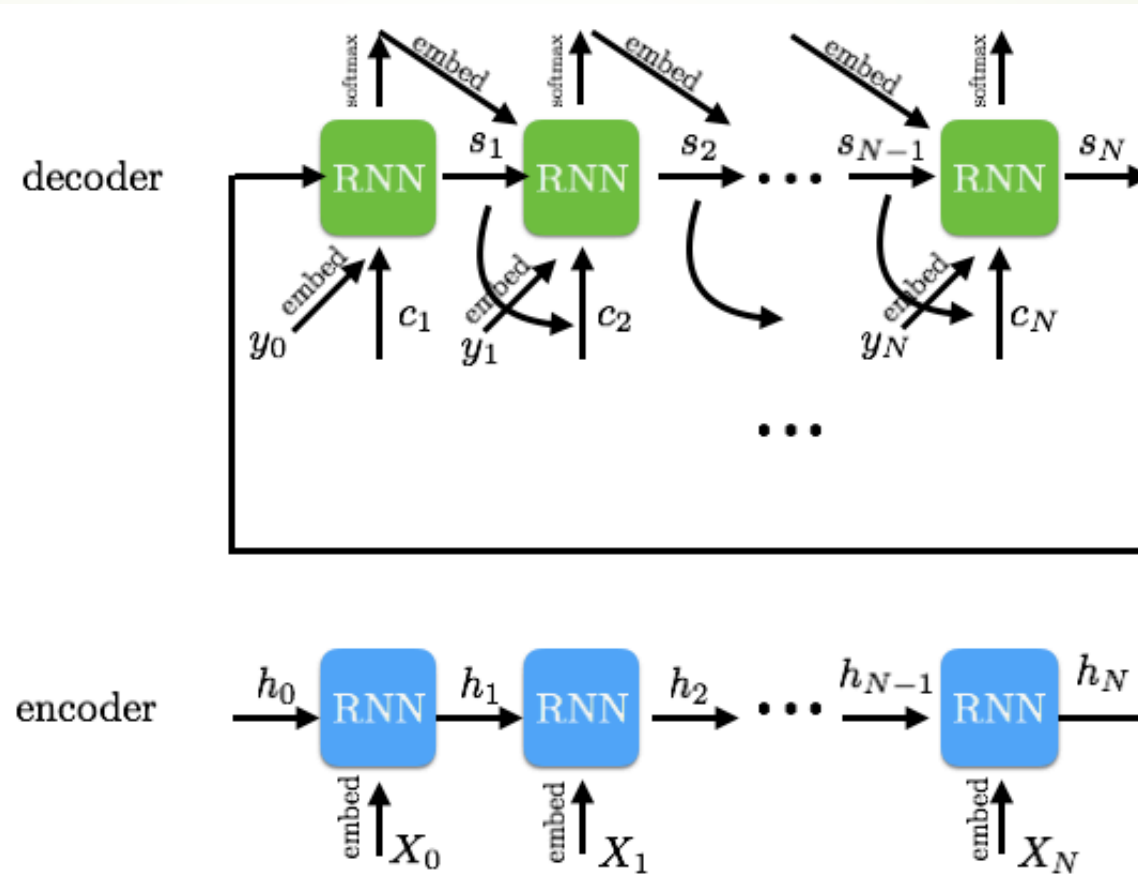
["Hello", ",", "World", "!"] → [1,4,5,6,7,2,0,...,0]

# Data Preprocess: Generate Files

- Each file contains 1024 articles or abstracts
- 183 files in training set
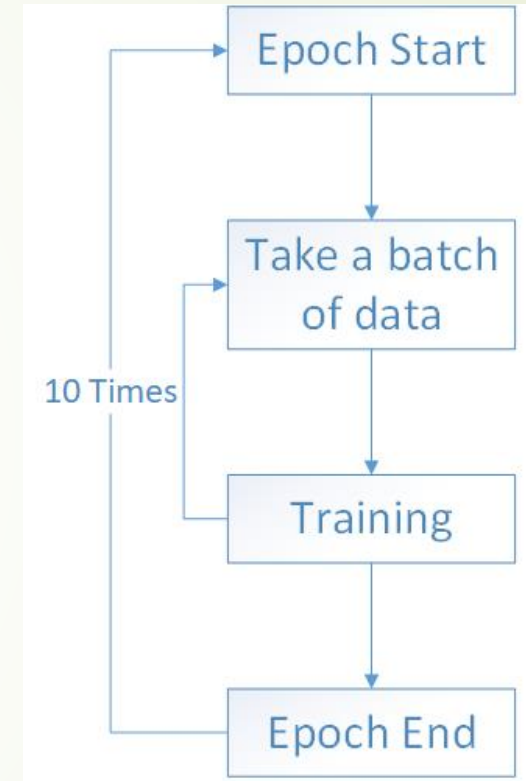- 61 files in testing and validation sets.

# Encoder/Decoder

# Training

- 10 epochs
- 4 documents each batch
- Trained 30 hours
- Loss decreases from 12.8 to 6.4

# Training

- Improve
  - Data parallel
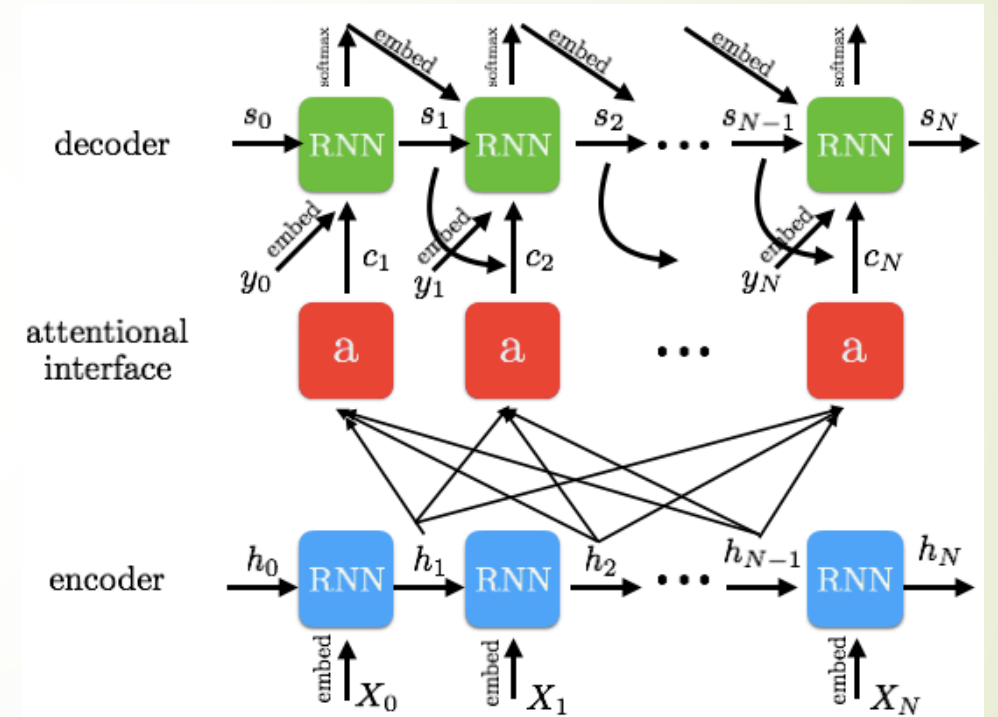  - 100 epochs
  - Flexible learning rate

# Testing

- Input the article and get predicted abstract
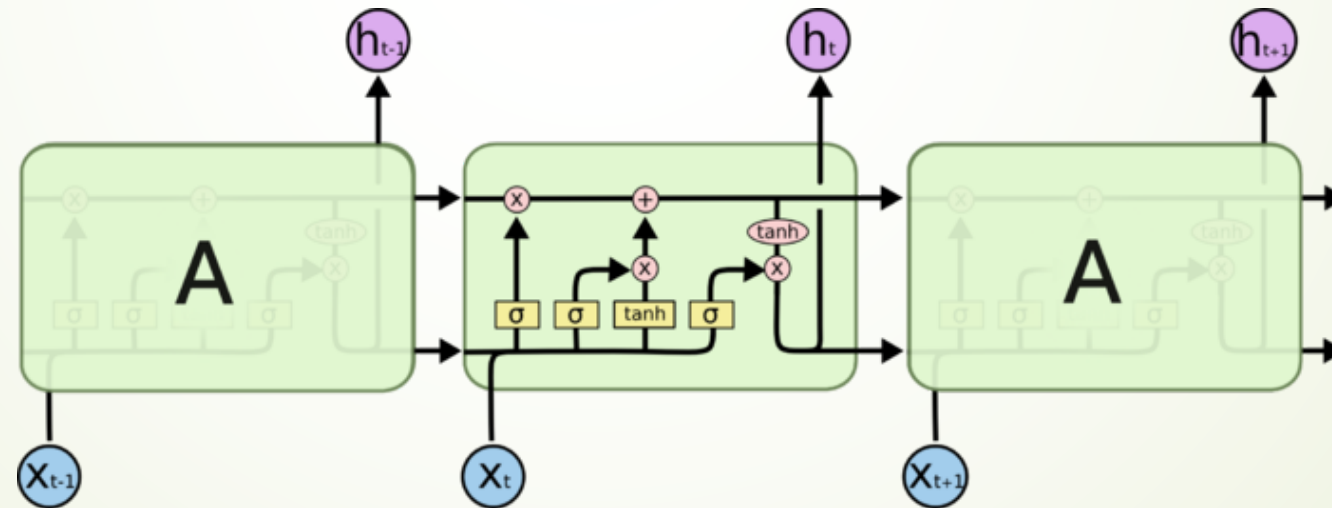- Evaluate by PyRouge

# Attention Model

- Based on the attention of human
- Take part of hidden value from encoder

# Long-Short Term Memory Model

- Based on human memory mechanism
- Encoder generate multiple hidden values
- Example: "The cloud is in the ___"

# Acknowledgements

- Client: Yufeng Ma

# References

- Encoder-decoder: https://theneuralperspective.com/2016/11/20/recurrent-neural-networks-rnn-part-3-encoder-decoder/

- Attention model: https://theneuralperspective.com/2016/11/20/recurrent-neural-network-rnn-part-4-attentional-interfaces/

- LSTM model: http://colah.github.io/posts/2015-08-Understanding-LSTMs/