

Deep learning for natural language processing

Advanced architectures

Benoit Favre <benoit.favre@univ-mrs.fr>

Aix-Marseille Université, LIF/CNRS

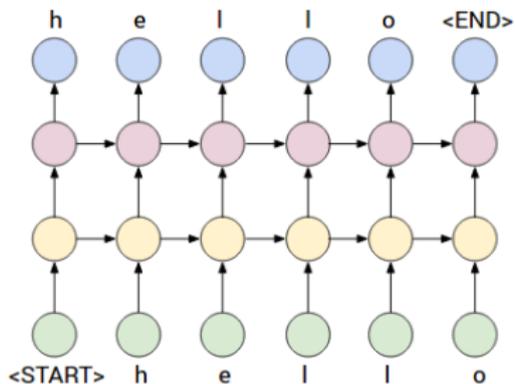
23 Feb 2017

Deep learning for Natural Language Processing

- Day 1
 - ▶ Class: intro to natural language processing
 - ▶ Class: quick primer on deep learning
 - ▶ Tutorial: neural networks with Keras
- Day 2
 - ▶ Class: word representations
 - ▶ Tutorial: word embeddings
- Day 3
 - ▶ Class: convolutional neural networks, recurrent neural networks
 - ▶ Tutorial: sentiment analysis
- Day 4
 - ▶ **Class: advanced neural network architectures**
 - ▶ Tutorial: language modeling
- Day 5
 - ▶ Tutorial: Image and text representations
 - ▶ Test

Stacked RNNs

- Increasing hidden state size is very expensive
 - ▶ U is of size ($hidden \times hidden$)
 - ▶ Can feed the output of a RNN to another RNN cell
 - ▶ → Multi-resolution analysis, better generalization

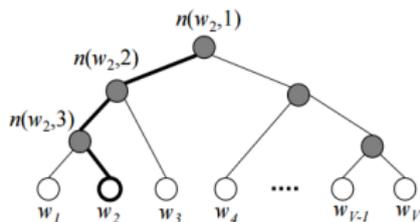


Source: <https://i.stack.imgur.com/usSPN.png>

- Necessary for large-scale language models

Softmax approximations

- When vocabulary is large (> 10000), the softmax layer gets too expensive
 - ▶ Store a $h \times |V|$ matrix in GPU memory
 - ▶ Training time gets very long
- Turn the problem to a sequence of decisions
 - ▶ Hierarchical softmax



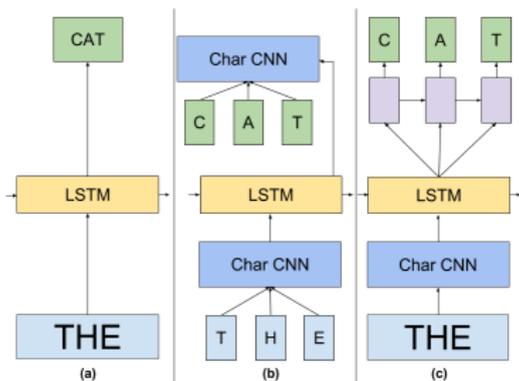
Source: <https://shuuki4.files.wordpress.com/2016/01/hsexample.png?w=1000>

- Turn the problem to a small set of binary decisions
 - ▶ Noise contrastive estimation, sampled softmax...
 - ▶ \rightarrow Pair target against a small set of randomly selected words
- More here: <http://sebastianruder.com/word-embeddings-softmax/>

Limits of language modeling

- Train a language model on the One Billion Word benchmark
 - ▶ “Exploring the Limits of Language Modeling”, Jozefowicz et al. 2016
 - ▶ 800k different words
 - ▶ Best model → 3 weeks on 32 GPU
 - ▶ PPL: perplexity evaluation metric (lower is better)

System	PPL
RNN-2048	68.3
Interpolated KN 5-GRAM	67.6
LSTM-512	32.2
2-layer LSTM-2048	30.6
Last row + CNN inputs	30.0
Last row + CNN softmax	39.8



Caption generation

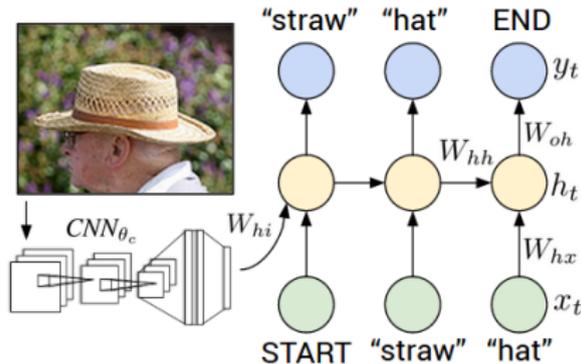
- Language model conditioned on an image
 - ▶ Generate image representation with CNN trained to recognize visual concepts
 - ▶ Stack image representation with language model input



people skiing on a snowy mountain



a woman playing tennis

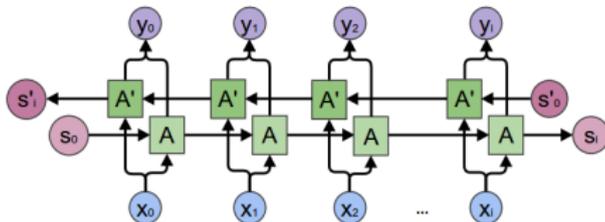


Source: <http://cs.stanford.edu/people/karpathy/rnn7.png>

- More here: <https://github.com/karpathy/neuraltalk2>

Bidirectional networks

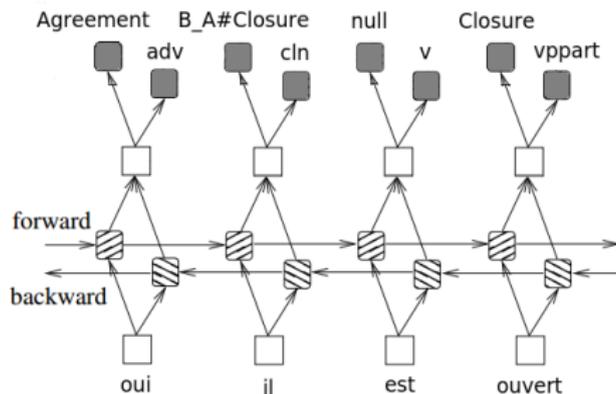
- RNN make predictions independent of the future observations
 - ▶ Need to look into the future
- Idea: concatenate the output of a forward and backward RNN
 - ▶ The decision can benefit from both past and future observations
 - ▶ Only applicable if we can wait for the future to happen



Source: <http://colah.github.io/posts/2015-09-NN-Types-FP/img/RNN-bidirectional.png>

Multi-task learning

- Can we build better representations by training the NN to predict different things?
 - ▶ Share the weights of lower system, diverge after representation layer
 - ▶ Also applies to feed forward neural networks
- Example: semantic tagging from words
 - ▶ Train system to predict low-level and high-level syntactic labels, as well as semantic labels
 - ▶ Need training data for each task
 - ▶ At test time only keep output of interest



Machine translation (the legacy approach)

Definitions

- *source* : text in the source language (ex: Chinese)
- *target* : text in the target language (ex: English)

Phrase-based statistical translation

- Decouple word translation and word ordering

$$P(\textit{target}|\textit{source}) = \frac{P(\textit{source}|\textit{target}) \times P(\textit{target})}{P(\textit{source})}$$

Model components

- $P(\textit{source}|\textit{target})$ = translation model
- $P(\textit{target})$ = language model
- $P(\textit{source})$ = ignored because constant for an input

Translation model

How to compute $P(\text{source}|\text{target}) = P(s_1, \dots, s_n | t_1, \dots, t_n)$?

$$P(s_1, \dots, s_n | t_1, \dots, t_n) = \frac{nb(s_1, \dots, s_n \rightarrow t_1, \dots, t_n)}{\sum_x nb(x \rightarrow t_1, \dots, t_n)}$$

- Piecewise translation

$$\begin{aligned} P(\text{I am your father} \rightarrow \text{Je suis ton père}) &= P(\text{I} \rightarrow \text{je}) \times P(\text{am} \rightarrow \text{suis}) \\ &\quad \times P(\text{your} \rightarrow \text{ton}) \\ &\quad \times P(\text{father} \rightarrow \text{père}) \end{aligned}$$

- To compute those probabilities
 - ▶ Need for alignment between source and target words

Alignements

I am your father
Je suis ton père

the boy **was looking** by the window
le garçon **regardait** par la fenêtre

He builds **houses**
Il construit **des maisons**

I am **not** like you
Je **ne** suis **pas** comme toi

It's raining **cats and dogs**
Il pleut **des cordes**

Have you done it yet ?
L'avez-vous déjà fait ?

They sell houses for a living
? Leur métier est de vendre des maisons

- Use bi-texts and alignment algorithm (such as Giza++)

Phrase table

	nous	ne	savons	pas	ce	qui	se	passe	.
we	■								
do		■							
not			■	■					
know			■						
what					■	■			
is							■		
happening								■	
.									■

	nous	ne	savons	pas	ce	qui	se	passe	.
we	■								
do		■							
not			■	■					
know			■						
what					■	■			
is							■		
happening								■	
.									■

	nous	ne	savons	pas	ce	qui	se	passe	.
we	■								
do		■							
not			■	■					
know			■						
what					■	■			
is							■		
happening								■	
.									■

"Phrase table"

we > nous
do not know > ne savons pas
what > ce qui
is happening > se passe
we do not know > nous ne savons pas
what is happening > ce qui se passe

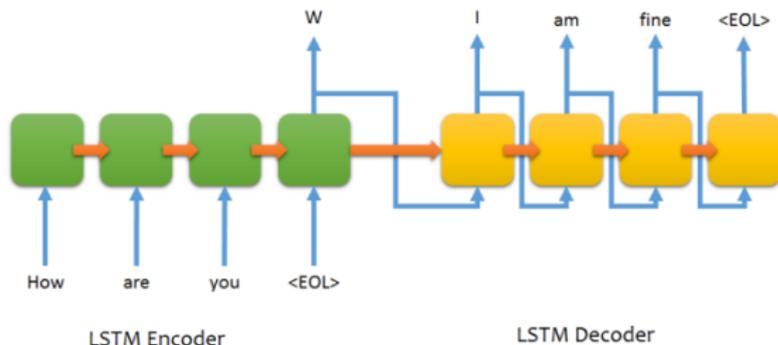
- Compute translation probability for all known phrases (an extension of n-gram language models)
 - ▶ Combine with LM and find best translation with decoding algorithm

Neural machine translation (NMT)

- Phrase-based translation
 - ▶ Same coverage problem as with word-ngrams
 - ▶ Alignment still wrong in 30% of cases
 - ▶ A lot of tricks to make it work
 - ▶ Researchers have progressively introduced NN
 - ★ Language model
 - ★ Phrase translation probability estimation
 - ▶ The google translate approach until mid-2016
- End-to-end approach to machine translation
 - ▶ Can we directly input source words and generate target words?

Encoder-decoder framework

- Generalisation of the conditioned language model
 - ▶ Build a representation, then generate sentence
 - ▶ Also called the seq2seq framework

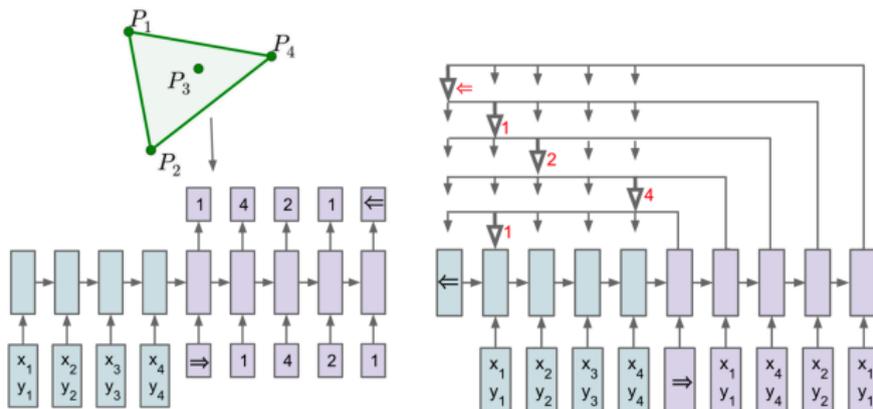


Source: <https://github.com/farizrahman4u/seq2seq>

- But still limited for translation
 - ▶ Bad for long sentences
 - ▶ How to account for unknown words?
 - ▶ How to make use of alignments?

Interlude: Pointer networks

- Decision is an offset in the input
 - ▶ Number of classes dependent on the length of the input
 - ▶ Decision depends on hidden state in input and hidden state in output
 - ▶ Can learn simple algorithms, such as finding the convex hull of a set of points



Source: <http://www.itdadao.com/articles/c19a1093068p0.html>

Oriol Vinyals, Meire Fortunato, Navdeep Jaitly, "Pointer Networks", arXiv:1506.03134

Attention mechanisms

- Loosely based on human visual attention mechanism
 - Let neural network focus on aspects of the input to make its decision
 - Learn what to attend based on what it has produced so far
 - More of a mechanism for memorizing the input

enc_j = encoder hidden state

dec_t = decoder hidden state

$$u_t^j = v^T \tanh(W_e enc_j + W_d dec_t)$$

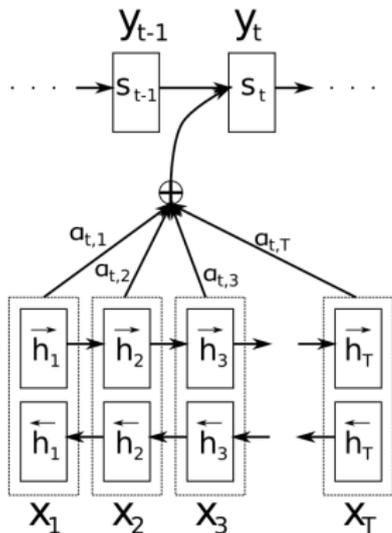
$$\forall j \in [1..n]$$

$$\alpha_t = \text{softmax}(u_t)$$

$$s_t = dec_t + \sum_j \alpha_t^j enc_j$$

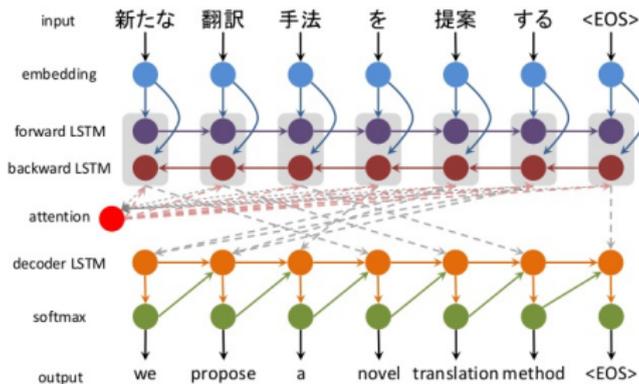
$$y_t = \text{softmax}(W_o s_t + b_o)$$

- New parameters: W_e, W_d, v



Machine translation with attention

Attention-based Neural Machine Translation



Source: <https://image.slidesharecdn.com/nmt-161019012948/95/attentionbased-nmt-description-4-638.jpg?cb=1476840773>

- Learns the word-to-word alignment

How to deal with unknown words

- If you don't have attention
 - ▶ Introduce *unk* symbols for low frequency words
 - ▶ Realign them to the input *a posteriori*
 - ▶ Use large translation dictionary or copy if proper name
- Use attention MT, extract α as alignment parameter
 - ▶ Then translate input word directly
- What about morphologically rich languages?
 - ▶ Reduce vocabulary size by translating word factors
 - ★ Byte pair encoding algorithm
 - ▶ Use word-level RNN to transliterate word

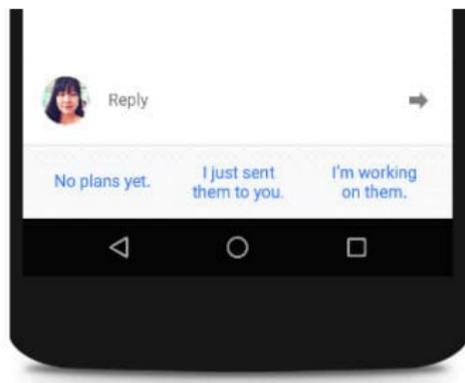
Zero-shot machine translation

- How to deal with the quadratic need for parallel data?
 - ▶ n languages $\rightarrow n^2$ pairs
 - ▶ So far, people have been using a pivot language ($x \rightarrow \text{english} \rightarrow y$)
- Parameter sharing across language pairs
 - ▶ Many to one \rightarrow share the target weights
 - ▶ One to many \rightarrow share the source weights
 - ▶ Many to many \rightarrow train single system for all pairs
- Zero-shot learning
 - ▶ Use token to identify target language (ex: <to-french>)
 - ▶ Let model learn to recognize source language
 - ▶ Can process pairs never seen in training!
 - ▶ The model learns the "interlingua"
 - ▶ Can also handle code switching

"Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation",
Johnson et al., arXiv:1611.04558

Conversation as translation

- Can we translate a question to its answer?
 - ▶ "Hello, how are you?" → "I am fine, thank you."
 - ▶ "What is the largest planet in the solar system?" → "It is Jupiter."
- "A Neural Conversational Model", Vinyals et al, 2015
 - ▶ Train a seq2seq model to generate the next turn in a dialog
 - ▶ Led to the "auto answer" feature in Google Inbox



Source: <http://cdn.ghacks.net/wp-content/uploads/2015/11/google-inbox-smart-reply.jpg>

What is a chatbot?

- Dialog system which can have an entertaining conversation
 - ▶ Chat-chat
 - ▶ Task oriented
- History
 - ▶ Eliza, virtual therapist
 - ★ <http://www.masswerk.at/elizabot/>
 - ▶ Mitsuku (best chatbot at Loebner price 2013/2016)
 - ★ <http://www.mitsuku.com/>
 - ▶ The Microsoft Tay fiasco
 - ★ Humans will always try to defeat an IA
 - ▶ A new industry hype
 - ★ Facebook, google...
- Question: can we spare dialog model engineering?
 - ▶ Train a model directly from conversation traces

Related work

- Models

- ▶ Generate next turn given previous turn with an encoder-decoder
 - ★ "A Neural Conversational Model" [Vynials et al. 2015]
- ▶ Add turn-level representations
 - ★ "Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models" [Serban et al., AAAI 2016]
- ▶ Add attention mechanism to the hierarchical model
 - ★ "Attention with Intention for a Neural Network Conversation Model" [Yao et al., SLUNIPS-2015]
- ▶ Chatbot as information retrieval
 - ★ "Improved Deep Learning Baselines for Ubuntu Corpus Dialogs" [Kadlec et al., SLUNIPS-2015]

- Dialog specifics

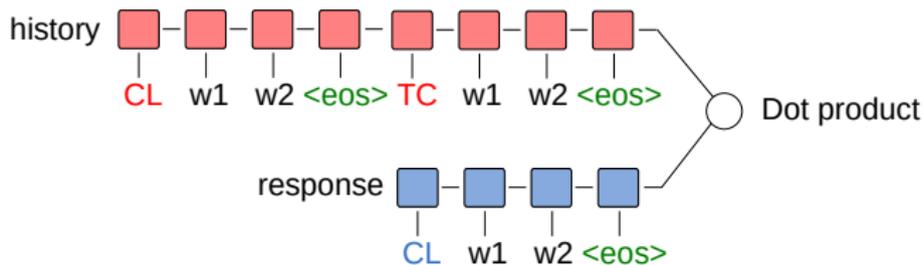
- ▶ Introduce long term reward
 - ★ "Deep Reinforcement Learning for Dialogue Generation", [Li et al., ACL 2016]
- ▶ How generate diverse responses?
 - ★ "A Diversity-Promoting Objective Function for Neural Conversation Models" [Li et al., NAACL 2016]
- ▶ Enforce consistency by explicitly modeling speakers
 - ★ "A Persona-Based Neural Conversation Model" [Li et al., ACL 2016]

- Evaluation: automatic metrics do not correlate with manual evaluation

- ▶ "How NOT To Evaluate Your Dialogue System" [Liu et al, EMNLP 2016]

Chatbot 2: bi-encoder

- Learn a model that gives the same representation to an answer and the context that led to it
 - ▶ Information retrieval which can retrieve the next turn given a history
 - ▶ Encode history with a first recurrent model
 - ▶ Encode next turn with a second recurrent model
 - ▶ Compute a similarity between those representations (dot product)
- Training objective
 - ▶ Make sure the correct association has a higher score than a randomly selected pair
- Problem: the cost of retrieving a turn
 - ▶ Everything can be precomputed, just the dot product remains
 - ▶ Many approaches for finding approximate nearest neighbors in a high dimensional space (ie. locality preserving hashing)

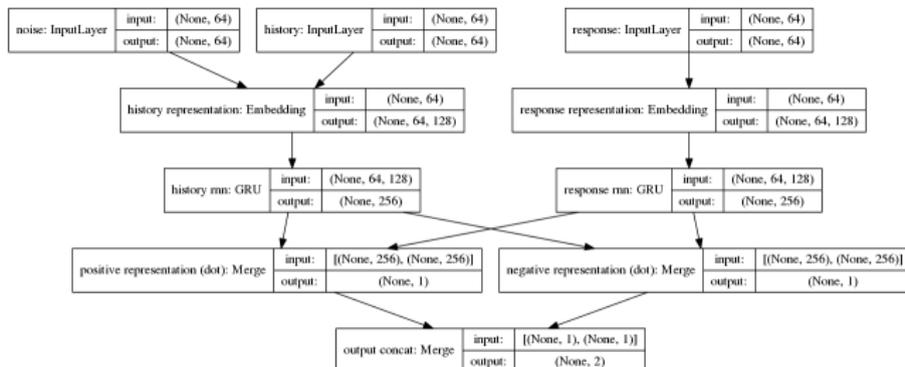


Bi-encoder training

- Maximize margin between the result of $h_i \cdot r_i$ and $n_i \cdot r_i$
 - ▶ h_i is the history
 - ▶ n_i is a random history
 - ▶ r_i is the response

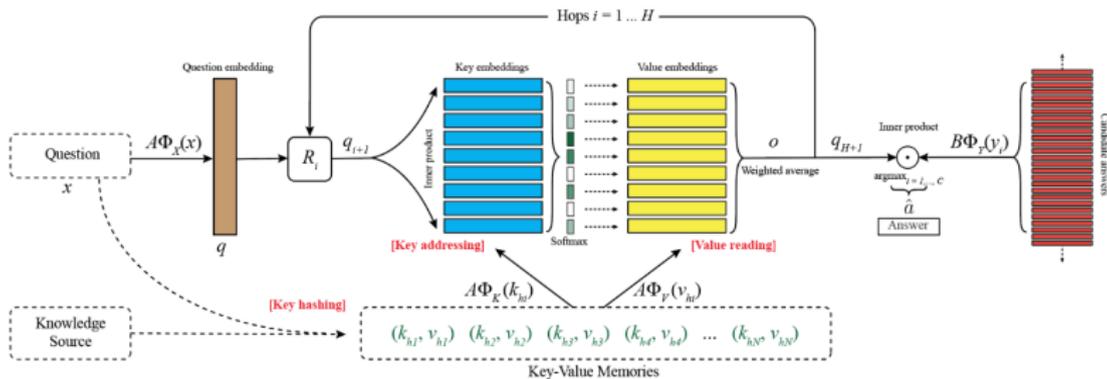
$$Loss = \frac{1}{n} \sum_i \max(0, 1 - h_i \cdot r_i + n_i \cdot r_i)$$

- Keras model



The future

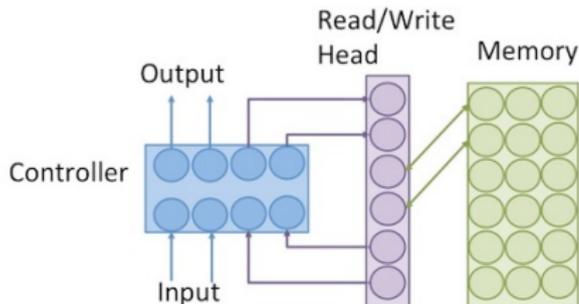
- Limitations of RNNs
 - ▶ Rewrite their memory at every time step
 - ▶ They have a fixed size memory
 - ▶ They need to reuse the same location in memory to perform the same action
- What if we had better memory devices
 - ▶ Static memory: Memory Networks (Weston et al., 2014)
 - ★ Memory containing representations (learned as part of the model)
 - ★ The model can do multiple passes over the memory to “deduce” its output



Source: <http://www.thespermwhale.com/jasveston/icml2016/mem1.png>

The future

- Dynamic memory: Neural Turing Machines
 - ▶ At each round
 - ★ Get memory read address from previous round
 - ★ Combine input, state and memory into new memory
 - ★ Generate memory read address for next round
 - ▶ Can learn basic algorithms
 - ★ Copy, sort...



Source: <http://lh3.googleusercontent.com/-Q0ZMLPrbLkU/ViucASG4HrI/AAAAAAAAABk4/-ZL4sny1-g0/s532-Ic42/nml1.jpeg>

Conclusion

- Add more prediction power to RNNs
 - ▶ Stacking
 - ▶ Bidirectional
 - ▶ Multitask
- Make better use of the input
 - ▶ Attention mechanisms
- Fancy applications
 - ▶ Machine translation
 - ▶ Caption generation
 - ▶ Chatbots

Remaining challenges

Deep learning for NLP

- Language independence
 - ▶ We still need training data in all languages
- Domain adaptation
 - ▶ Often, we have plenty of data where we don't need it, and none where we would need it
 - ▶ What if the test data does not follow the distribution of training data?
- Dealing with small datasets
 - ▶ Annotating complex phenomena is expensive

Deep learning

- Efficient training on CPU, mobile devices
 - ▶ Binary neural networks
- Training non differentiable systems
 - ▶ Reinforcement learning
- Reasoning, world knowledge...
 - ▶ AI, here we are