# Memory Networks
## Deep Learning Lecture 10

Samuel Cheng

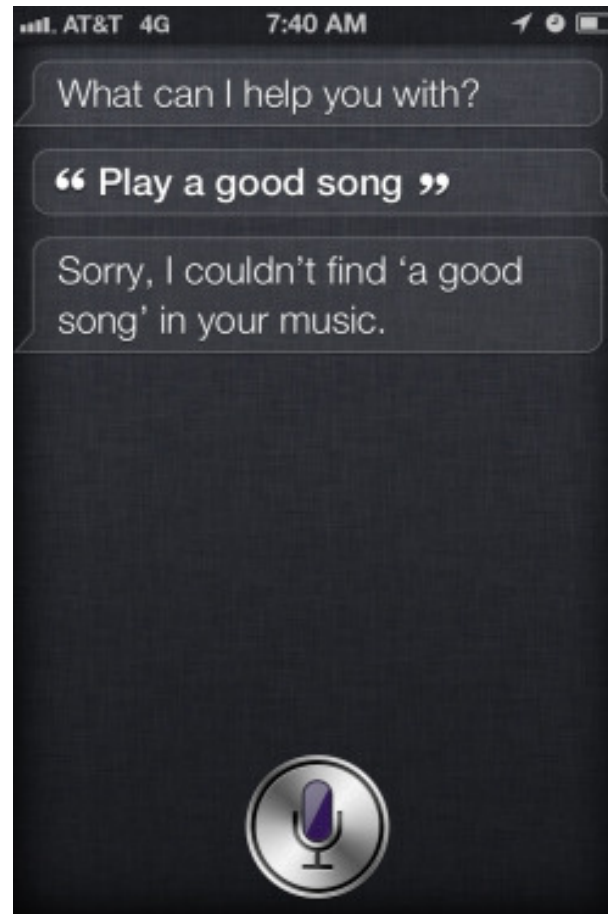**Slides credit to Jason Weston**

# Logistics

- Activity 2 was up. As usual, first successful submission will be awarded by 3% overall bonus
  - As winner of last activity, Siraj will be excluded for the competition

- Muhanad and Siraj will present Caffe today. Please vote promptly

- Next class will be back to the normal time (Friday 3:30 pm)

# Review and Overview

- We discussed the sequence-to-sequence models last week. And we also looked into two examples
  - Neural machine translation
  - Chatbot

- We will talk about memory networks today. As an example, we will also look into Q&A systems and dialogue agents

# Intelligent Conversational Agents

# End-to-End Dialog Agents

While we already have some kind of useful dialog agents (*Google Now, Cortana, Siri, .. ?*) a true dialog agent should:

- Be able to **combine all its knowledge (both common facts and user-specific info)** *to fulfill complex tasks*.

- **Handle long open-ended conversations** *involving effectively tracking many latent variables*.

- Be able to **learn** (new tasks) via conversation.

*Opinion from Facebook AI Research (FAIR):* **Machine Learning End-to-End systems** is the way forward in the long-run.

# Memory Networks

- Class of models that combine large memory with learning component that can read and write to it.

- Incorporates **reasoning** with **attention** over **memory** (RAM).

- Most ML has limited memory which is more-or-less all that's needed for "low level" tasks e.g. object detection.

**Motivation:** long-term memory is required to read a story and then e.g. answer questions about it.

Similarly, it's also required for **dialog**: to remember previous dialog (short- and long-term), and respond.

# Memory Networks

| | | |
|---|---|---|
| Long-Term Memories | $h_i$ | Shaolin Soccer directed_by Stephen Chow |
| | | Shaolin Soccer written_by Stephen Chow |
| | | Shaolin Soccer starred_actors Stephen Chow |
| | | Shaolin Soccer release_year 2001 |
| | | Shaolin Soccer has_genre comedy |
| | | Shaolin Soccer has_tags martial arts, kung fu soccer, stephen chow |
| | | Kung Fu Hustle directed_by Stephen Chow |
| | | Kung Fu Hustle written_by Stephen Chow |
| | | Kung Fu Hustle starred_actors Stephen Chow |
| | | Kung Fu Hustle has_genre comedy action |
| | | Kung Fu Hustle has_imdb_votes famous |
| | | Kung Fu Hustle has_tags comedy, action, martial arts, kung fu, china, soccer, hong kong, stephen chow |
| | | The God of Cookery directed_by Stephen Chow |
| | | The God of Cookery written_by Stephen Chow |
| | | The God of Cookery starred_actors Stephen Chow |
| | | The God of Cookery has_tags hong kong Stephen Chow |
| | | From Beijing with Love directed_by Stephen Chow |
| | | From Beijing with Love written_by Stephen Chow |
| | | From Beijing with Love starred_actors Stephen Chow, Anita Yuen |
| | | ... <and more> ... |
| Short-Term Memories | $c_1^u$ $c_1^r$ | 1) I'm looking a fun comedy to watch tonight, any ideas? 2) Have you seen Shaolin Soccer? That was zany and great.. really funny but in a whacky way. |
| Input | $c_2^u$ | 3) Yes! Shaolin Soccer and Kung Fu Hustle are so good I really need to find some more Stephen Chow films I feel like there is more awesomeness out there that I haven't discovered yet ... |
| Output | $y$ | 4) God of Cookery is pretty great, one of his mid 90's hong kong martial art comedies. |

# What is a Memory Network?
## *Original paper description of class of models*

MemNNs have four component networks (which may or may not have shared parameters):

- **I:** (input feature map) convert incoming data to the internal feature representation.

- **G:** (generalization) update memories given new input.

- **O:** produce new output (in feature representation space) given the memories.

- **R:** (response) convert output O into response seen by the outside world.
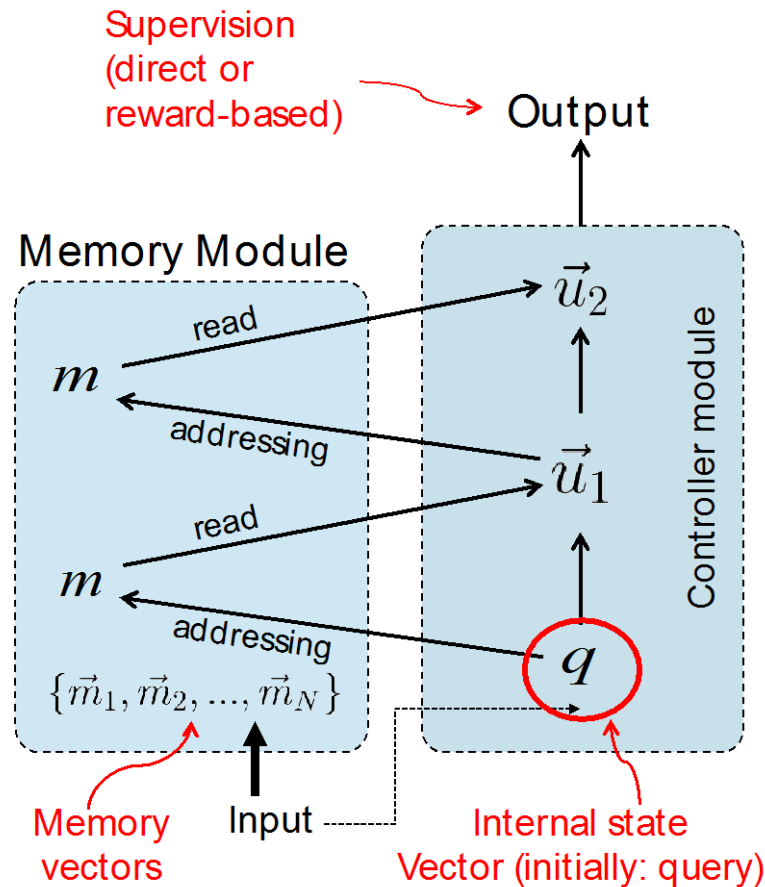
# Some Memory Network-related Publications

- J. Weston, S. Chopra, A. Bordes. Memory Networks. ICLR 2015 (and arXiv:1410.3916).

- S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. End-To-End Memory Networks. NIPS 2015 (and arXiv:1503.08895).

- J. Weston, A. Bordes, S. Chopra, A. M. Rush, B. van Merriënboer, A. Joulin, T. Mikolov. Towards AI-Complete Question Answering: A Set of Prerequisite Toy Tasks. arXiv:1502.05698.

- A. Bordes, N. Usunier, S. Chopra, J. Weston. Large-scale Simple Question Answering with Memory Networks. arXiv:1506.02075.

- J. Dodge, A. Gane, X. Zhang, A. Bordes, S. Chopra, A. Miller, A. Szlam, J. Weston. Evaluating Prerequisite Qualities for Learning End-to-End Dialog Systems. arXiv:1511.06931.

- F. Hill, A. Bordes, S. Chopra, J. Weston. The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations. arXiv:1511.02301.

- J. Weston. Dialog-based Language Learning. arXiv:1604.06045.

- A. Bordes, Jason Weston. Learning End-to-End Goal-Oriented Dialog. arXiv:1605.07683.

# Memory Network Models
## *schematic diagram*



[Figure by Saina Sukhbaatar]

# Variants of the class…

Some options and extensions:

- **Representation of inputs and memories could use all kinds of encodings:** bag of words, RNN style reading at word or character level, etc.

- **Different possibilities for output module:** e.g. multi-class classifier or uses an RNN to output sentences.

- **If the memory is huge** (e.g. Wikipedia) we need to organize the memories. Solution: hash the memories to store in buckets (topics). Then, memory addressing and reading doesn't operate on *all* memories.

- **If the memory is full**, there could be a way of removing one it thinks is most useless; *i.e.* it ``forgets'' somehow. That would require a scoring function of the utility of each memory..

# Task (1) Factoid QA with Single Supporting Fact ("where is actor")

(Very Simple) Toy reading comprehension task:

John was in the bedroom.
Bob was in the office.
John went to kitchen. ← SUPPORTING FACT
Bob travelled back home.
Where is John? A:kitchen

# (2) Factoid QA with Two Supporting Facts ("where is actor+object")

A harder (toy) task is to answer questions where two supporting statements have to be chained to answer the question:

John is in the playground.
Bob is in the office.
John picked up the football.
Bob went to the kitchen.
Where is the football?  A:playground
Where was Bob before the kitchen? A:office

# (2) Factoid QA with Two Supporting Facts ("where is actor+object")

A harder (toy) task is to answer questions where two supporting statements have to be chained to answer the question:

John is in the playground.  ←——— SUPPORTING FACT
Bob is in the office.
John picked up the football.  ←——— SUPPORTING FACT
Bob went to the kitchen.
Where is the football?  A:playground
Where was Bob before the kitchen? A:office

To answer the first question *Where is the football?* both *John picked up the football* and *John is in the playground* are supporting facts.

.

# The First **MemNN** Implemention

- I (input): converts to bag-of-word-embeddings $x$.

- G (generalization): stores $x$ in next available slot $m_N$.

- O (output): Loops over all memories k=1 or 2 times:
  - 1st loop max: finds best match $m_i$ with $x$.
  - 2nd loop max: finds best match $m_J$ with $(x, m_i)$.
  - The output $o$ is represented with $(x, m_i, m_J)$.

- R (response): ranks all words in the dictionary given $o$ and returns best single word. *(OR: use a full RNN here)*

# Matching function

- For a given Q, we want a good match to the relevant memory slot(s) containing the answer, e.g.:

Match(Where is the football ?,  John picked up the football)

- We use a $q^TU^TUd$ embedding model with word embedding features:

  - *LHS features:*  Q:Where Q:is Q:the Q:football Q:?

  - *RHS features:*  D:John D:picked D:up D:the D:football

**The parameters U are trained with a margin ranking loss: supporting facts should score higher than non-supporting facts.**

# Matching function: 2$^{nd}$ hop

- On the 2$^{nd}$ hop we match question & 1$^{st}$ hop to new fact:

Match( [Where is the football ?, John picked up the football],
        John is in the playground)

- We use the same q$^T$U$^T$Ud embedding model:
  - *LHS features:*  Q:Where Q:is Q:the Q:football Q:? Q2: John Q2:picked Q2:up Q2:the Q2:football
  - *RHS features:*  D:John D:is D:in D:the D:playground

# Objective function

Minimize:

$$\sum_{\bar{f} \neq \mathbf{m}_{o_1}} \max(0, \gamma - s_O(x, \mathbf{m}_{o_1}) + s_O(x, \bar{f})) +$$

$$\sum_{\bar{f}' \neq \mathbf{m}_{o_2}} \max(0, \gamma - s_O([x, \mathbf{m}_{o_1}], \mathbf{m}_{o_2}) + s_O([x, \mathbf{m}_{o_1}], \bar{f}')) +$$

$$\sum_{\bar{r} \neq r} \max(0, \gamma - s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], r) + s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], \bar{r}))$$

Where: $S_O$ is the matching function for the Output component.
$S_R$ is the matching function for the Response component.
x is the input question.
$m_{O1}$ is the first true supporting memory (fact).
$m_{O2}$ is the first second supporting memory (fact).
r is the response
True facts and responses $m_{O1}$, $m_{O2}$ and r should have higher scores than all other facts and responses by a given margin.

# Comparing triples

- We also need time information for the bAbI tasks. We tried adding absolute time as a feature: it works, but it is better to compare using triplets

- Seems to work better if we compare triples:

- Match(Q,D,D') returns < 0 if D is better than D'

returns > 0 if D' is better than D

*We can loop through memories, keep best $m_i$ at each step.*

# Comparing triples: Objective and Inference

$$\sum_{\bar{f} \neq \mathbf{m}_{o_1}} \max(0, \gamma - s_{O_t}(x, \mathbf{m}_{o_1}, \bar{f})) + \sum_{\bar{f} \neq \mathbf{m}_{o_1}} \max(0, \gamma + s_{O_t}(x, \bar{f}, \mathbf{m}_{o_1})) +$$

$$\sum_{\bar{f}' \neq \mathbf{m}_{o_2}} \max(0, \gamma - s_{O_t}([x, \mathbf{m}_{o_1}], \mathbf{m}_{o_2}, \bar{f}')) + \sum_{\bar{f}' \neq \mathbf{m}_{o_2}} \max(0, \gamma + s_{O_t}([x, \mathbf{m}_{o_1}], \bar{f}', \mathbf{m}_{o_2}) +$$

$$\sum_{\bar{r} \neq r} \max(0, \gamma - s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], r) + s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], \bar{r}))$$

$$o_1 = O_t(x, \mathbf{m}) \qquad o_2 = O_t([x, m_{o_1}], \mathbf{m})$$

---

**Algorithm 1** $O_t$ replacement to arg max when using write time features

---

**function** $O_t(q, \mathbf{m})$
    $t \leftarrow 1$
    **for** $i = 2, \ldots, N$ **do**
        **if** $s_{O_t}(q, \mathbf{m}_i, \mathbf{m}_t) > 0$ **then**
            $t \leftarrow i$
        **end if**
    **end for**
    **return** $t$
**end function**

---

# bAbI Experiment 1

- 10k sentences. (Actor: only ask questions about actors.)
- Difficulty: how many sentences in the past when entity mentioned.
- Fully supervised (supporting sentences are labeled).
- Compare RNN (no supervision)
  and MemNN hops $k = 1$ or $2$, & with/without time features.

| Method | Difficulty 1 | | Difficulty 5 | |
| --- | --- | --- | --- | --- |
| | actor | actor+object | actor | actor+object |
| RNN | 0% | 42% | 71% | 83% |
| MemNN $k = 1$ | 10% | 81% | 54% | 79% |
| MemNN $k = 1$ (+time) | 0% | 27% | 0% | 27% |
| MemNN $k = 2$ (+time) | 0% | 0.05% | 0% | 0.6% |

*Difficulty 5 -- Max mem. sz. required: 65   Average mem. sz. required: 9*

# bAbI Experiment 1

- Example test story + predictions:

Antoine went to the kitchen. Antoine got the milk. Antoine travelled to the office. Antoine dropped the milk. Sumit picked up the football. Antoine went to the bathroom. Sumit moved to the kitchen.

- *where is the milk now?* A: office

- *where is the football?* A: kitchen

- *where is Antoine ?* A: bathroom

- *where is Sumit ?* A: kitchen

- *where was Antoine before the bathroom*? A: office

# Larger QA: Reverb Dataset in (Fader et al., 13)

- 14M statements, stored as (subject, relation, object) triples. Triples are REVERB extractions mined from ClueWeb09.

- Statements cover diverse topics:
  - (milne, authored, winnie-the-pooh)
  - (sheep, be-afraid-of, wolf),   *etc...*

- Weakly labeled QA pairs and 35M paraphrased questions from WikiAnswers:
  - ``Who wrote the Winnie the Pooh books?
  - ``Who is poohs creator?''

# Results: QA on Reverb data from (Fader et al.)

- 14M statements stored in the memNN memory.
- *k=1 loops* MemNN, *128-dim* embedding.
- R response simply outputs top scoring statement.
- Time features are not necessary, hence not used.
- We also tried adding bag of words (BoW) features.

| Method | F1 |
|---|---|
| (Fader et al., 2013) | 0.54 |
| (Bordes et al., 2014) | 0.73 |
| MemNN | 0.72 |
| MemNN (with BoW features) | 0.82 |

# Fast QA on Reverb data

Scoring all 14M candidates in the memory is slow.

We consider speedups using hashing in S and O as mentioned earlier:

- Hashing via words (*essentially: inverted index*)

- Hashing via k-means in embedding space (k=1000)

| Method | Embedding | Embed+BoW | candidates |
|---|---|---|---|
| MemNN (no hashing) | 0.72 | 0.82 | 14M |
| MemNN (word hash) | 0.63 | 0.68 | 13k (1000x) |
| MemNN (clust hash) | 0.71 | 0.80 | 177k (80x) |

Embed+BoW: $q^T U^T U d + \lambda\, q^T d$

# A MemNN multitasked on bAbI data and Reverb QA data

**The "story" told to the model *after* training:**

*Antoine went to the kitchen. Antoine picked up the milk. Antoine travelled to the office.*

**MemNN's answers to some questions:**

- *Where is the milk?* A: office

- *Where was Antoine before the office?* A: kitchen

- *Where does milk come from?* A: milk come from cow

- *What is a cow a type of?* A: cow be female of cattle

- *Where are cattle found?* A: cattle farm become widespread in brazil

- *What does milk taste like?* A: milk taste like milk

- *What does milk go well with?* A: milk go with coffee

# What was next for MemNNs?

- Make the language much harder: coreference, conjunctions, negations, etc. etc – *will it work?*

- MemNNs that reason with *more than* 2 supporting memories.

- End-to-end? (doesn't need supporting facts for training)

- More useful applications on real datasets.

- Dialog: Ask questions? Say statements?

- *Do MemNN ideas extend to other ML tasks and model variants, .e.g. visual QA, perform actions…? [A: yes!].*

# bAbI tasks: what reasoning tasks would we like models to work on?

- We define 20 tasks (generated by the simulation) that we can test new models on. (See: http://fb.ai/babi)

- The idea is they are a bit like software tests: each task checks if an ML system has a certain skill.

- We would like each "skill" we check to be a natural task for humans w.r.t. text understanding & reasoning, humans should be able to get 100%.

J. Weston, A. Bordes, S. Chopra, T. Mikolov. Towards AI-Complete Question Answering: A Set of Prerequisite Toy Tasks. arXiv:1502.05698.

# Simulation commands

```
go <place>

get <object>

get <object1> from <object2>

put <object1> in/on <object2>

give <object> to <person>

drop <object>

look

inventory

examine <object>
```

*+ 2 commands for "gods" (superusers):*

```
create <object>

set <obj1> <relation> <obj2>
```

# Example

**Command format**

```
jason go kitchen

jason get milk

jason go office

jason drop milk

jason go bathroom

where is milk ?    A: office

where is jason? A: bathroom
```

**Story**

Jason went to the kitchen.

Jason picked up the milk.

Jason travelled to the office.

Jason left the milk there.

Jason went to the bathroom.

Where is the milk now? A: office

Where is Jason? A: bathroom

# Task (1) Factoid QA with Single Supporting Fact ("where is actor")

Our first task consists of questions where a single supporting fact, previously given, provides the answer.

We test simplest case of this, by asking for the location of a person.

A small sample of the task is thus:

John is in the playground. ← SUPPORTING FACT
Bob is in the office.
Where is John? A:playground

We could use supporting facts for supervision at training time, but are not known at test time (we call this "strong supervision"). However weak supervision is much better!!

# (2) Factoid QA with Two Supporting Facts ("where is actor+object")

A harder task is to answer questions where two supporting statements have to be chained to answer the question:

John is in the playground.  ←——————  SUPPORTING FACT
Bob is in the office.
John picked up the football.  ←——————  SUPPORTING FACT
Bob went to the kitchen.
Where is the football?  A:playground

To answer the question *Where is the football?* both *John picked up the football* and *John is in the playground* are supporting facts.

.

# (3) Factoid QA with Three Supporting Facts

Similarly, one can make a task with three supporting facts:

John picked up the apple.
John went to the office.
John went to the kitchen.
John dropped the apple.
Where was the apple before the kitchen? A:office

The first three statements are all required to answer this.

# (4) Two Argument Relations: Subject vs. Object

To answer questions the ability to differentiate and recognize subjects and objects is crucial.

We consider the extreme case: sentences feature re-ordered words:

The office is north of the bedroom.
The bedroom is north of the bathroom.
What is north of the bedroom? A:office
What is the bedroom north of? A:bathroom

Note that the two questions above have exactly the same words, but in a different order, and different answers.

So a bag-of-words will not work.

# (6) Yes/No Questions

- This task tests, in the simplest case possible (with a single supporting fact) the ability of a model to answer true/false type questions:

John is in the playground.
Daniel picks up the milk.
Is John in the classroom? A:no
Does Daniel have the milk? A:yes

# (7) Counting

Tests ability to count sets:

> Daniel picked up the football.
> Daniel dropped the football.
> Daniel got the milk.
> Daniel took the apple.
> How many objects is Daniel holding? A:two

# (8) Lists/Sets

- Tests ability to produce lists/sets:

> Daniel picks up the football.
> Daniel drops the newspaper.
> Daniel picks up the milk.
> What is Daniel holding? A:milk,football

# (11) Basic Coreference (nearest referent)

Daniel was in the kitchen.
Then he went to the studio.
Sandra was in the office.
Where is Daniel? A:studio

# (13) Compound Coreference

Daniel and Sandra journeyed to the office.
Then they went to the garden.
Sandra and John travelled to the kitchen.
After that they moved to the hallway.
Where is Daniel? A:garden

In linguistics, coreference, sometimes written co-reference, occurs when two or more expressions in a text refer to the same person or thing; they have the same referent, e.g. Bill said he would come; the proper noun Bill and the pronoun he refer to the same person, namely to Bill.

# (14) Time manipulation

- While our tasks so far have included time implicitly in the *order* of the statements, this task tests understanding the use of time expressions within the statements:

> In the afternoon Julie went to the park.
> Yesterday Julie was at school.
> Julie went to the cinema this evening.
> Where did Julie go after the park? A:cinema

**Much harder difficulty:** adapt a real time expression labeling dataset into a question answer format, e.g. Uzzaman et al., '12.

# (15) Basic Deduction

- This task tests basic deduction via inheritance of properties:

Sheep are afraid of wolves.
Cats are afraid of dogs.
Mice are afraid of cats.
Gertrude is a sheep.
What is Gertrude afraid of? A:wolves

Deduction should prove difficult for MemNNs because it effectively involves search, although our setup might be simple enough for it.

# (17) Positional Reasoning

- This task tests spatial reasoning, one of many components of the classical SHRDLU system:

The triangle is to the right of the blue square.
The red square is on top of the blue square.
The red sphere is to the right of the blue square.
Is the red sphere to the right of the blue square? A:yes
Is the red square to the left of the triangle? A:yes

# (18) Reasoning about size

- This tasks requires reasoning about relative size of objects and is inspired by the commonsense reasoning examples in the Winograd schema challenge:

The football fits in the suitcase.
The suitcase fits in the cupboard.
The box of chocolates is smaller than the football.
Will the box of chocolates fit in the suitcase? A:yes

Tasks 3 (three supporting facts) and 6 (Yes/No) are prerequisites.

# (19) Path Finding

- In this task the goal is to find the path between locations:

> The kitchen is north of the hallway.
> The den is east of the hallway.
> How do you go from den to kitchen?  A:west,north

This is going to prove difficult for MemNNs because it effectively involves search.

# End-to-end Memory Network (MemN2N)

- New end-to-end (MemN2N) model (Sukhbaatar '15):
  - Reads from memory with **soft attention**
  - Performs **multiple lookups** (hops) on memory
  - End-to-end training with **backpropagation**
  - Only need supervision on the final output

- Compared to original "Memory Networks" (MemNN), which had:
  - Hard attention
  - requires explicit supervision of attention during training
  - Only feasible for simple tasks

# MemN2N schematic diagram

# Memory Module

# Question & Answering

Answer: kitchen

## Memory Module



Weighted Sum $\rightarrow$ $0.1\vec{m}_1 + 0.7\vec{m}_2 + 0.2\vec{m}_3$ $\rightarrow$ $\vec{u}_2$

$\{0.1, 0.7, 0.2\}$

Dot product + softmax $\leftarrow$ $\vec{u}_1$

$\{\vec{m}_1, \vec{m}_2, \vec{m}_3\}$

Controller

1: Sam moved to garden

2: Sam went to kitchen

3: Sam drops apple there

Where is Sam?

Input story

Question

# Memory Vectors

E.g.) constructing memory vectors with Bag-of-Words (BoW)

1. Embed each word

2. Sum embedding vectors

"Sam drops apple" $\rightarrow \underbrace{\vec{v}_{\text{Sam}} + \vec{v}_{\text{drops}} + \vec{v}_{\text{apple}}}_{\text{Embedding Vectors}}$

# Positional Encoding of Words

**Representation of inputs and memories could use all kinds of encodings:** bag of words, RNN style reading at word or character level, etc.

**We also built a positional encoding variant:** Words are represented by vectors as before. But instead of a bag, position is modeled by a multiplicative term on each word vector with weights depending on the position in the sentence.

# Training on 1k stories

| TASK | Weakly supervised | | | Supervised Supp. Facts | |
|---|---|---|---|---|---|
| | N-grams | LSTMs | MemN2N | Memory Networks | StructSVM +coref+srl |
| T1. Single supporting fact | 36 | 50 | PASS | PASS | PASS |
| T2. Two supporting facts | 2 | 20 | 87 | PASS | 74 |
| T3. Three supporting facts | 7 | 20 | 60 | PASS | 17 |
| T4. Two arguments relations | 50 | 61 | PASS | PASS | PASS |
| T5. Three arguments relations | 20 | 70 | 87 | PASS | 83 |
| T6. Yes/no questions | 49 | 48 | 92 | PASS | PASS |
| T7. Counting | 52 | 49 | 83 | 85 | 69 |
| T8. Sets | 40 | 45 | 90 | 91 | 70 |
| T9. Simple negation | 62 | 64 | 87 | PASS | PASS |
| T10. Indefinite knowledge | 45 | 44 | 85 | PASS | PASS |
| T11. Basic coreference | 29 | 72 | PASS | PASS | PASS |
| T12. Conjunction | 9 | 74 | PASS | PASS | PASS |
| T13. Compound coreference | 26 | PASS | PASS | PASS | PASS |
| T14. Time reasoning | 19 | 27 | PASS | PASS | PASS |
| T15. Basic deduction | 20 | 21 | PASS | PASS | PASS |
| T16. Basic induction | 43 | 23 | PASS | PASS | 24 |
| T17. Positional reasoning | 46 | 51 | 49 | 65 | 61 |
| T18. Size reasoning | 52 | 52 | 89 | PASS | 62 |
| T19. Path finding | 0 | 8 | 7 | 36 | 49 |
| T20. Agent's motivation | 76 | 91 | PASS | PASS | PASS |

# Attention during memory lookups

## Samples from toy QA tasks

| Story (1: 1 supporting fact) | Support | Hop 1 | Hop 2 | Hop 3 |
|---|---|---|---|---|
| Daniel went to the bathroom. | | 0.00 | 0.00 | 0.03 |
| Mary travelled to the hallway. | | 0.00 | 0.00 | 0.00 |
| John went to the bedroom. | | 0.37 | 0.02 | 0.00 |
| John travelled to the bathroom. | yes | 0.60 | 0.98 | 0.96 |
| Mary went to the office. | | 0.01 | 0.00 | 0.00 |
| Where is John?   Answer: bathroom   Prediction: bathroom | | | | |

| Story (2: 2 supporting facts) | Support | Hop 1 | Hop 2 | Hop 3 |
|---|---|---|---|---|
| John dropped the milk. | | 0.06 | 0.00 | 0.00 |
| John took the milk there. | yes | 0.88 | 1.00 | 0.00 |
| Sandra went back to the bathroom. | | 0.00 | 0.00 | 0.00 |
| John moved to the hallway. | yes | 0.00 | 0.00 | 1.00 |
| Mary went back to the bedroom. | | 0.00 | 0.00 | 0.00 |
| Where is the milk?   Answer: hallway   Prediction: hallway | | | | |

| Story (16: basic induction) | Support | Hop 1 | Hop 2 | Hop 3 |
|---|---|---|---|---|
| Brian is a frog. | yes | 0.00 | 0.98 | 0.00 |
| Lily is gray. | | 0.07 | 0.00 | 0.00 |
| Brian is yellow. | yes | 0.07 | 0.00 | 1.00 |
| Julius is green. | | 0.06 | 0.00 | 0.00 |
| Greg is a frog. | yes | 0.76 | 0.02 | 0.00 |
| What color is Greg?   Answer: yellow   Prediction: yellow | | | | |

| Story (18: size reasoning) | Support | Hop 1 | Hop 2 | Hop 3 |
|---|---|---|---|---|
| The suitcase is bigger than the chest. | yes | 0.00 | 0.88 | 0.00 |
| The box is bigger than the chocolate. | | 0.04 | 0.05 | 0.10 |
| The chest is bigger than the chocolate. | yes | 0.17 | 0.07 | 0.90 |
| The chest fits inside the container. | | 0.00 | 0.00 | 0.00 |
| The chest fits inside the box. | | 0.00 | 0.00 | 0.00 |
| Does the suitcase fit in the chocolate?   Answer: no   Prediction: no | | | | |

20 bAbI Tasks

| | Test Acc | Failed tasks |
|---|---|---|
| MemNN | 93.3% | 4 |
| LSTM | 49% | 20 |
| MemN2N 1 hop | 74.82% | 17 |
| 2 hops | 84.4% | 11 |
| 3 hops | 87.6.% | 11 |

# So we still fail on some tasks….

*.. and we could also make more tasks that we fail on!*

Our hope is that a feedback loop of:

1. Developing **tasks** that **break models**, and

2. Developing **models** that can **solve tasks**

   … leads in a **fruitful** research direction….

# How about on real data?

- Toy AI tasks are important for developing innovative methods.

- But they do not give all the answers.

- How do these models work on real data?
  - Classic Language Modeling (Penn TreeBank, Text8)
  - Story understanding (Children's Book Test, News articles)
  - Open Question Answering (WebQuestions, WikiQA)
  - Goal-Oriented Dialog and Chit-Chat (Movie Dialog, Ubuntu)

# Language Modeling

The goal is to predict the next word in a text sequence given the previous words. Results on the Penn Treebank and Text8 (Wikipedia-based) corpora.

|               | Penn Tree | Text8 |
|---------------|-----------|-------|
| RNN           | 129       | 184   |
| LSTM          | 115       | 154   |
| MemN2N 2 hops | 121       | 187   |
| 5 hops        | 118       | 154   |
| 7 hops        | 111       | 147   |

Test perplexity

Hops vs. Attention:
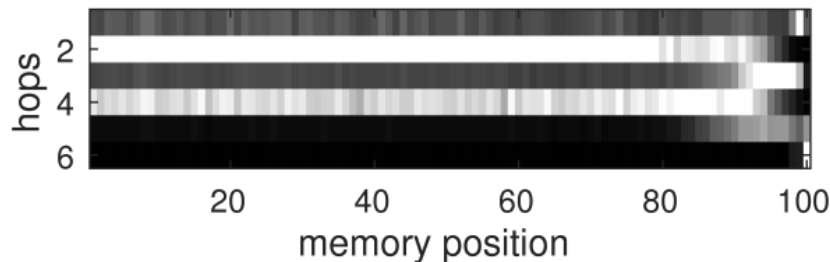Average over (PTB)          Average over (Text8)

# Language Modeling

The goal is to predict the next word in a text sequence given the previous words. Results on the Penn Treebank and Text8 (Wikipedia-based) corpora.

|  | Penn Tree | Text8 |
|---|---|---|
| RNN | 129 | 184 |
| LSTM | 115 | 154 |
| MemN2N 2 hops | 121 | 187 |
| 5 hops | 118 | 154 |
| 7 hops | 111 | 147 |

Test perplexity

**MemNNs are in the same ballpark as LSTMs.**

**Hypothesis: many words (e.g. syntax words) don't actually need really long term context, and so memNNs don't help there.**

**Maybe MemNNs could eventually help more on things like nouns/entities?**

# Children books understanding

New dataset based on 118 children books from project Gutenberg

growing increasingly alarmed at the likelihood of their neocolony falling to English-speaking rebels. In mid-June, just as my hotel was being evacuated, the French announced plans to send a peace-keeping mission to the western part of Rwanda for "humanitarian" reasons. This gave the génocidaires the chance to look like victims instead of aggressors, and they started to pack up and leave for the protected area that became known as "the Turquoise Zone."

RTLM radio then performed its final disservice to the nation by scaring the living daylights out of the people remaining in Rwanda, a considerable number of whom had just spent two months murdering their neighbors and chasing the less compliant ones through swamps. The radio told them that the RPF would kill any Hutus they found in their path and encouraged all its listeners to pack up their belongings and head toward Goma or the western part of the country and the borders of the ... Republic of Congo (what used to be called Zaire), where the French soldiers awaited. Nearly 1.7 million people heeded the call. Entire hills and cities mobilized into caravans: men carrying sacks of bananas, some with bloody machetes in their belt loops; women with baskets of grain on their heads; children hugging photo albums to their chests. They ... corpses piled at the side of the road and the smouldering cooking fires in front of looted houses. I am sorry to say that the dire predictions of the radio were not rooted in fantasy, as the rebels did conduct crimes against humanity in revenge for the genocide and to make people fear them. In any case, what was left of Rwanda emptied out within days.

The U.N. Security Council, so ineffective in the face of the genocide, lent its sponsorship to the camps the French set up to protect the "refugees." The main place of comfort to the killers was at a town called Goma, just over the border into the Democratic ...

**Context**

**Question**

canoes and t hellish lands equipped pa jets, tents, we pathetic UN height in Ap shelter some
Many of parently thes attack the ne the Interahan the camps, p keep filling th camp so thei faithful. It w comfort was
In a surp suaded to act ten adminis for the camp ple who occ initiative to c over into Uga times what it which would corpses.
On July 4, RPF captured conquered a were knocked were empty al

1 " phebe beckoned to him ; i saw her , " cried rose , staring hard at the door .

2 " is it more presents coming ? "

3 asked jamie , just as his brother re-appeared , looking more excited than ever .

4 " yes ; a present for mother , and here it is ! "

5 roared archie , flinging wide the door to let in a tall man , who cried out , " where 's my little woman ?

6 the first kiss for her , then the rest may come on as fast as they like . "

7 before the words were out of his mouth , mrs. jessie was half-hidden under his rough great-coat , and four boys were prancing about him clamouring for their turn .

8 of course , there was a joyful tumult for a time , during which rose slipped into the window recess and watched what went on , as if it were a chapter in a christmas story .

9 it was good to see bluff uncle jem look proudly at his tall son , and fondly hug the little ones .

10 it was better still to see him shake his brothers ' hands as if he would never leave off , and kiss all the sisters in a way that made even solemn aunt myra brighten up for a minute .

11 but it was best of all to see him finally established in grandfather 's chair , with his " little woman " beside him , his three youngest boys in his lap , and _____ hovering over him like a large-sized cherub .

faith | brothers | rose | **archie** | rest | mouth | way | mother | sisters | george

# MemNNs for story understanding



**MemNN memory**

| | |
|---|---|
| | . . . . . |
| $m_0$ | NULL |
| | `` Why , what are YOUR shoes done with ? ' said the Gryphon |
| | 'I mean , what makes them so shiny ? |
| | Alice looked down at them , and considered a little before she gave her answer . . |
| $m_n$ | ..... |

**Memory reads and stores story**

**Story**

`` Why , what are YOUR shoes done with ? ' said the Gryphon . 'I mean , what makes them so shiny ? ' Alice looked down at them , and considered a little before she gave her answer . `` They 're done with blacking , I believe . ' `` Boots and shoes under the sea , ' the Gryphon went on in a deep voice , `` are done with a whiting . Now you know . ' `` And what are they made of ? ' Alice asked in a tone of great curiosity .. '

Cands: Gryphon | Alice | King | Queen | ...

`` Soles and eels , of course , ' the _____ replied rather impatiently : `` any shrimp could have told you that

Gryphon

# Results on Children's Book Test

| Methods | Named Entities | Common Nouns | Verbs | Prepositions |
|---|---|---|---|---|
| Humans (query)[*] | 0.520 | 0.644 | 0.716 | 0.676 |
| Humans (context+query)[*] | *0.816* | *0.816* | *0.828* | 0.708 |
| Maximum frequency (corpus) | 0.120 | 0.158 | 0.373 | 0.315 |
| Maximum frequency (context) | 0.335 | 0.281 | 0.285 | 0.275 |
| Sliding window | 0.168 | 0.196 | 0.182 | 0.101 |
| Word distance model | 0.398 | 0.364 | 0.380 | 0.237 |
| Kneser-Ney language model | 0.390 | 0.544 | 0.778 | 0.768 |
| Kneser-Ney language model + cache | 0.439 | 0.577 | 0.772 | 0.679 |
| Embedding Model (context+query) | 0.253 | 0.259 | 0.421 | 0.315 |
| Embedding Model (query) | 0.351 | 0.400 | 0.614 | 0.535 |
| Embedding Model (window) | 0.362 | 0.415 | 0.637 | 0.589 |
| Embedding Model (window+position) | 0.402 | 0.506 | 0.736 | 0.670 |
| LSTMs (query) | 0.408 | 0.541 | 0.813 | 0.802 |
| LSTMs (context+query) | 0.418 | 0.560 | **0.818** | 0.791 |
| Contextual LSTMs (window context) | 0.436 | 0.582 | 0.805 | **0.806** |
| MemNNs (lexical memory) | 0.431 | 0.562 | 0.798 | 0.764 |
| MemNNs (window memory) | 0.493 | 0.554 | 0.692 | 0.674 |
| MemNNs (sentential memory + PE) | 0.318 | 0.305 | 0.502 | 0.326 |
| MemNNs (window memory + self-sup.) | **0.666** | **0.630** | 0.690 | 0.703 |

# Question Answering on New's Articles

We evaluate our models on the data from:
**"Teaching Machines to Read and Comprehend"**
Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt,
Will Kay, Mustafa Suleyman, Phil Blunsom

| Original Version | Anonymised Version |
|---|---|
| **Context** | |
| The BBC producer allegedly struck by Jeremy Clarkson will not press charges against the "Top Gear" host, his lawyer said Friday. Clarkson, who hosted one of the most-watched television shows in the world, was dropped by the BBC Wednesday after an internal investigation by the British broadcaster found he had subjected producer Oisin Tymon "to an unprovoked physical and verbal attack." ... | the *ent381* producer allegedly struck by *ent212* will not press charges against the " *ent153* " host , his lawyer said friday . *ent212* , who hosted one of the most - watched television shows in the world , was dropped by the *ent381* wednesday after an internal investigation by the *ent180* broadcaster found he had subjected producer *ent193* " to an unprovoked physical and verbal attack . " ... |
| **Query** | |
| Producer **X** will not press charges against Jeremy Clarkson, his lawyer says. | producer **X** will not press charges against *ent212* , his lawyer says . |
| **Answer** | |
| Oisin Tymon | *ent193* |

# Results on CNN QA dataset

| METHODS | VALIDATION | TEST |
|---|---|---|
| MAXIMUM FREQUENCY (ARTICLE)[*] | 0.305 | 0.332 |
| SLIDING WINDOW | 0.005 | 0.006 |
| WORD DISTANCE MODEL[*] | 0.505 | 0.509 |
| DEEP LSTMS (ARTICLE+QUERY)[*] | 0.550 | 0.570 |
| CONTEXTUAL LSTMS ("ATTENTIVE READER")[*] | 0.616 | 0.630 |
| CONTEXTUAL LSTMS ("IMPATIENT READER")[*] | 0.618 | 0.638 |
| MEMNNS (WINDOW MEMORY) | 0.580 | 0.606 |
| MEMNNS (WINDOW MEMORY + SELF-SUP.) | 0.634 | 0.668 |
| MEMNNS (WINDOW MEMORY + ENSEMBLE) | 0.612 | 0.638 |
| MEMNNS (WINDOW MEMORY + SELF-SUP. + ENSEMBLE) | 0.649 | 0.684 |
| MEMNNS (WINDOW + SELF-SUP. + ENSEMBLE + EXCLUD. COOCURRENCES) | **0.662** | **0.694** |

Table 3: **Results on CNN QA.** [*]Results taken from Hermann et al. (2015).

# Large Scale QA

**MemNN memory**

**Memory *reads and stores* Freebase**

**Freebase**

22 M facts
5 M entities

......

```
Gollum character created by JRR Tolkien
JRR Tolkien place of birth Bloemfontein
Bloemfontein contained by South Africa
The Hobbit directed by Peter Jackson
Facebook Inc founded by Mark Zuckerberg
```

**Read Module *looks for* 1 *sup. fact* *among a subset: uses hashing/ string matching for fast lookup.***

Who created Gollum from The Hobbit?  |  Gollum character created by JRR_Tolkien

**R Module *returns the object***

JRR_Tolkien

# WebQuestions & SimpleQuestions

- Decent results on WebQuestions, a popular QA task:

| | |
|---|---|
| Random guess | 1.9 |
| (Bordes et al., 2014b) | 29.7 |
| (Berant et al., 2013) | 31.3 |
| (Bordes et al., 2014a) | 39.2 |
| (Berant and Liang, 2014) | 39.9 |
| (Yang et al., 2014) | 41.3 |
| MemNN | 42.2 |

A. Bordes, N. Usunier, S. Chopra J.Weston. Large-scale Simple Question Answering with Memory Networks. arXiv:1506.02075.

- However now beaten by many results, especially (Yih et al. ACL '15) that achieves **52.5**! Several hand engineered features are used in that case. Note WebQuestions is very small (4k train+valid).

# Recent Work: New Models for QA on documents Miller et al. Key-Value Memory Networks for Directly Reading Documents. arXiv:1606.03126.

# Recent Work: New Models for QA on documents Miller et al. Key-Value Memory Networks for Directly Reading Documents. arXiv:1606.03126.

**WikiQA Results**

| Method | MAP | MRR |
|---|---|---|
| Word Cnt | 0.4891 | 0.4924 |
| Wgt Word Cnt | 0.5099 | 0.5132 |
| 2-gram CNN (Yang *et al.*, 2015) | 0.6520 | 0.6652 |
| AP-CNN (Santos *et al.*, 2016) | 0.6886 | 0.6957 |
| Attentive LSTM (Miao *et al.*, 2015) | 0.6886 | 0.7069 |
| Attentive CNN (Yin and Schütze, 2015) | 0.6921 | 0.7108 |
| L.D.C. (Wang *et al.*, 2016) | 0.7058 | 0.7226 |
| Memory Network | 0.5170 | 0.5236 |
| **Key-Value Memory Network** | **0.7069** | **0.7265** |

Question

What year was the movie Blade Runner released?

Knowledge Source

**Wikipedia Entry: Blade Runner**

Blade Runner ▮▮▮▮▮▮ dystopian science fiction film ▮▮▮▮▮▮ Harrison Ford, Rutger Hauer, Sean Young, and Edward James Olmos. The screenplay, ▮▮▮▮▮▮ is a modified film adaptation of the 1968 novel "Do Androids Dream of Electric Sheep?" by Philip K. Dick. The film depicts a dystopian Los Angeles in November 2019 in which genetically engineered replicants, which are visually indistinguishable from adult humans, are manufactured by the powerful Tyrell Corporation as

Candidate answers

Tron

1982

police

Tom Cruise

...

ican neo-noir / 1982
ican neo-noir / Blade Runner
scott and starring / R. Scott
scott and starring / Blade Runner
ancher and D. P. / H. Fancher
ancher and D. P. / Blade Runner

# Recent Work: Combines QA with Dialog Tasks
## Dodge et al. "Evaluating Prerequisite Qualities for Learning End-to-End Dialog Systems." ICLR '16

**(Dialog 1) QA:** *facts about movies*

Sample input contexts and target replies (in red) from Dialog Task 1:

What movies are about open source?  Revolution OS
Ruggero Raimondi appears in which movies? Carmen
What movies did Darren McGavin star in?  Billy Madison, The Night
Stalker, Mrs. Pollifax-Spy, The Challenge
Can you name a film directed by Stuart Ortiz? Grave Encounters
Who directed the film White Elephant?  Pablo Trapero
What is the genre of the film Dial M for Murder?  Thriller, Crime
What language is Whity in?  German

**(Dialog 2) Recs:** *movie recommendations*

Sample input contexts and target replies (in red) from Dialog Task 2:

Schindler's List, The Fugitive, Apocalypse Now, Pulp Fiction, and
The Godfather are films I really liked. Can you suggest a film?
The Hunt for Red October

Some movies I like are Heat, Kids, Fight Club, Shaun of the Dead,
The Avengers, Skyfall, and Jurassic Park. Can you suggest
something else I might like?  Ocean's Eleven

**(Dialog 3) QA+Recs:** *combination dialog*

Sample input contexts and target replies (in red) from Dialog Task 3:

I loved Billy Madison, Blades of Glory, Bio-Dome, Clue, and Happy
Gilmore. I'm looking for a Music movie.  School of Rock
What else is that about?  Music, Musical, Jack Black, school,
teacher, Richard Linklater, rock, guitar
I like rock and roll movies more. Do you know anything else?
Little Richard

**(Dialog 4) Reddit:** *real dialog*

Sample input contexts and target replies (in red) from Dialog Task 4:

I think the Terminator movies really suck, I mean the first one was
kinda ok, but after that they got really cheesy. Even the second one
which people somehow think is great. And after that...
forgeddabotit.
C'mon the second one was still pretty cool.. Arny was still so
badass, as was Sararah Connor's character.. and the way they
blended real action and effects was perhaps the last of its kind...

# (Dialog 1) QA: *facts about movies*

Sample input contexts and target replies (in red) from Dialog Task 1:

What movies are about open source?  Revolution OS
Ruggero Raimondi appears in which movies? Carmen
What movies did Darren McGavin star in?  Billy Madison, The Night
Stalker, Mrs. Pollifax-Spy, The Challenge
Can you name a film directed by Stuart Ortiz? Grave Encounters
Who directed the film White Elephant?  Pablo Trapero
What is the genre of the film Dial M for Murder?   Thriller, Crime
What language is Whity in?  German

# (Dialog 2) Recs: *movie recommendations*

Sample input contexts and target replies (in red) from Dialog Task 2:

Schindler's List, The Fugitive, Apocalypse Now, Pulp Fiction, and The Godfather are films I really liked. Can you suggest a film?      The Hunt for Red October

Some movies I like are Heat, Kids, Fight Club, Shaun of the Dead, The Avengers, Skyfall, and Jurassic Park. Can you suggest something else I might like?  Ocean's Eleven

# (Dialog 3) QA+Recs: *combination dialog*

Sample input contexts and target replies (in red) from Dialog Task 3:

I loved Billy Madison, Blades of Glory, Bio-Dome, Clue, and Happy Gilmore. I'm looking for a Music movie.   School of Rock
What else is that about?    Music, Musical, Jack Black, school, teacher, Richard Linklater, rock, guitar
I like rock and roll movies more. Do you know anything else?
Little Richard

# (Dialog 4) Reddit: *real dialog*

Sample input contexts and target replies (in red) from Dialog Task 4:

I think the Terminator movies really suck, I mean the first one was kinda ok, but after that they got really cheesy. Even the second one which people somehow think is great. And after that... forgeddabotit.
C'mon the second one was still pretty cool.. Arny was still so badass, as was Sararah Connor's character.. and the way they blended real action and effects was perhaps the last of its kind...

# Memory Network: *example*

| | | |
|---|---|---|
| Memories | $h_i$ | Shaolin Soccer written_by Stephen Chow |
| | | Shaolin Soccer starred_actors Stephen Chow |
| | | Shaolin Soccer release_year 2001 |
| | | Shaolin Soccer has_genre comedy |
| | | Shaolin Soccer has_tags martial arts, kung fu soccer, stephen chow |
| | | Kung Fu Hustle directed_by Stephen Chow |
| | | Kung Fu Hustle written_by Stephen Chow |
| | | Kung Fu Hustle starred_actors Stephen Chow |
| | | Kung Fu Hustle has_genre comedy action |
| | | Kung Fu Hustle has_imdb_votes famous |
| | | Kung Fu Hustle has_tags comedy, action, martial arts, kung fu, china, soccer, hong kong, stephen chow |
| | | The God of Cookery directed_by Stephen Chow |
| | | The God of Cookery written_by Stephen Chow |
| | | The God of Cookery starred_actors Stephen Chow |
| | | The God of Cookery has_tags hong kong Stephen Chow |
| | | From Beijing with Love directed_by Stephen Chow |
| | | From Beijing with Love written_by Stephen Chow |
| | | From Beijing with Love starred_actors Stephen Chow, Anita Yuen |
| | | ...<and more> ... |
| Short-Term Memories | $c_1^u$ $c_1^r$ | 1) I'm looking a fun comedy to watch tonight, any ideas? <br> 2) Have you seen Shaolin Soccer? That was zany and great.. really funny but in a whacky way. |
| Input | $c_2^u$ | 3) Yes! Shaolin Soccer and Kung Fu Hustle are so good I really need to find some more Stephen Chow films I feel like there is more awesomeness out there that I haven't discovered yet ... |

# Results

| Methods | QA Task (hits@1) | Recs Task (hits@100) | QA+Recs Task (hits@10) | Reddit Task (hits@10) |
|---|---|---|---|---|
| QA System (Bordes et al., 2014) | 90.7 | N/A | N/A | N/A |
| SVD | N/A | 19.2 | N/A | N/A |
| IR | N/A | N/A | N/A | 23.7 |
| LSTM | 6.5 | 27.1 | 19.9 | 11.8 |
| Supervised Embeddings | 50.9 | 29.2 | 65.9 | 27.6 |
| MemN2N | 79.3 | 28.6 | 81.7 | 29.2 |
| Joint Supervised Embeddings | 43.6 | 28.1 | 58.9 | 14.5 |
| Joint MemN2N | 83.5 | 26.5 | 78.9 | 26.6 |

# Ubuntu Data

Dialog dataset: Ubuntu IRC channel logs, users ask questions about issues they are having with Ubuntu and get answers by other users. (Lowe et al., '15)

| METHODS | | VALIDATION (HITS@1) | TEST (HITS@1) |
|---|---|---|---|
| IR[†] | | N/A | 48.81 |
| RNN[†] | | N/A | 37.91 |
| LSTM[†] | | N/A | 55.22 |
| MEMN2N | 1-HOP | 57.23 | 56.25 |
| MEMN2N | 2-HOPS | 64.28 | 63.51 |
| MEMN2N | 3-HOPS | 64.31 | 63.72 |
| MEMN2N | 4-HOPS | 64.01 | 62.82 |

Table 7: **Ubuntu Dialog Corpus results.** The evaluation is retrieval-based, similar to that of Reddit (Task 4). For each dialog, the correct answer is mixed among 10 random candidates; Hits@1 (in %) are reported. Methods with [†] have been ran by Lowe et al. (2015).

Best results currently reported:
Sentence Pair Scoring: Towards Unified Framework for Text Comprehension
Petr Baudiš, Jan Pichl, Tomáš Vyskočil, Jan Šedivý
RNN-CNN combo model: 67.2

# FAIR: paper / data / code

- Papers:
  - bAbI tasks: arxiv.org/abs/1502.05698
  - Memory Networks: http://arxiv.org/abs/1410.3916
  - End-to-end Memory Networks: http://arxiv.org/abs/1503.08895
  - Large-scale QA with MemNNs: http://arxiv.org/abs/1506.02075
  - Reading Children's Books: http://arxiv.org/abs/1511.02301
  - Evaluating End-To-End Dialog:  http://arxiv.org/abs/1511.06931
  - Dialog-based Language Learning: http://arxiv.org/abs/1604.06045

- Data:
  - bAbI tasks: fb.ai/babi
  - SimpleQuestions dataset (100k questions): fb.ai/babi
  - Children's Book Test dataset: fb.ai/babi
  - Movie Dialog Dataest: fb.ai/babi

- Code:
  - Memory Networks: https://github.com/facebook/MemNN
  - Simulation tasks generator: https://github.com/facebook/bAbI-tasks

# RAM Issues

- **How to decide what to write and what not to write in the memory?**

- **How to represent knowledge to be stored in memories?**

- **Types of memory (arrays, stacks, or stored within weights of model), when they should be used, and how can they be learnt?**

- **How to do fast retrieval of relevant knowledge from memories when the scale is huge?**

- **How to build hierarchical memories, e.g. multiscale attention?**

- **How to build hierarchical reasoning, e.g. composition of functions?**

- **How to incorporate forgetting/compression of information?**

- **How to evaluate reasoning models? Are artificial tasks a good way? Where do they break down and real tasks are needed?**

- **Can we draw inspiration from how animal or human memories work?**

# Presentation continues next week

| Date | Student | Package |
|---|---|---|
| 3/3 | Aakash<br>Soubhi | Tensorflow<br>Tensorflow |
| 3/10 | Ahmad A<br>Tamer | Theano<br>Theano |
| 3/23 | Ahmad M<br>Obada | Keras<br>Keras |
| 3/30 | Muhanad<br>Siraj | Caffe<br>Caffe |
| 4/7 | **Dong**<br>**Varun** | Torch<br>Lasagne |
| 4/14 | Naim | MatConvNet |

# Thanks!