

Deep Reasoning

2016-03-16

Taehoon Kim

carpedm20@gmail.com

References

1. **[Sukhbaatar, 2015]** Sukhbaatar, Szlam, Weston, Fergus. *“End-To-End Memory Networks”* Advances in Neural Information Processing Systems. 2015.
2. **[Hill, 2015]** Hill, Bordes, Chopra, Weston. *“The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations”* arXiv preprint arXiv:1511.02301 (2015).
3. **[Kumar, 2015]** Kumar, Irsoy, Ondruska, Iyyer, Bradbury, Gulrajani, Zhong, Paulus, Socher. *“Ask Me Anything: Dynamic Memory Networks for Natural Language Processing”* arXiv preprint arXiv:1511.06038 (2015).
4. **[Xiong, 2016]** Xiong, Merity, Socher. *“Dynamic Memory Networks for Visual and Textual Question Answering”* arXiv preprint arXiv:1603.01417 (2016).
5. **[Hermann, 2015]** Hermann, Kočiský, Grefenstette, Espeholt, Will Kay, Suleyman, Blunsom. *“Teaching Machines to Read and Comprehend”* arXiv preprint arXiv:arXiv:1506.03340 (2015).
6. **[Miao, 2015]** Miao, Lei Yu, Blunsom. *“Neural Variational Inference for Text Processing”* arXiv preprint arXiv:1511.06038 (2015).
7. **[Kingma, 2013]** Kingma, Diederik P., and Max Welling. *“Auto-encoding variational bayes”* arXiv preprint arXiv:1312.6114 (2013).
8. **[Sohn, 2015]** Sohn, Kihyuk, Honglak Lee, and Xinchen Yan. *“Learning Structured Output Representation using Deep Conditional Generative Models.”* Advances in Neural Information Processing Systems. 2015.

Models

Not Considered transitive inference

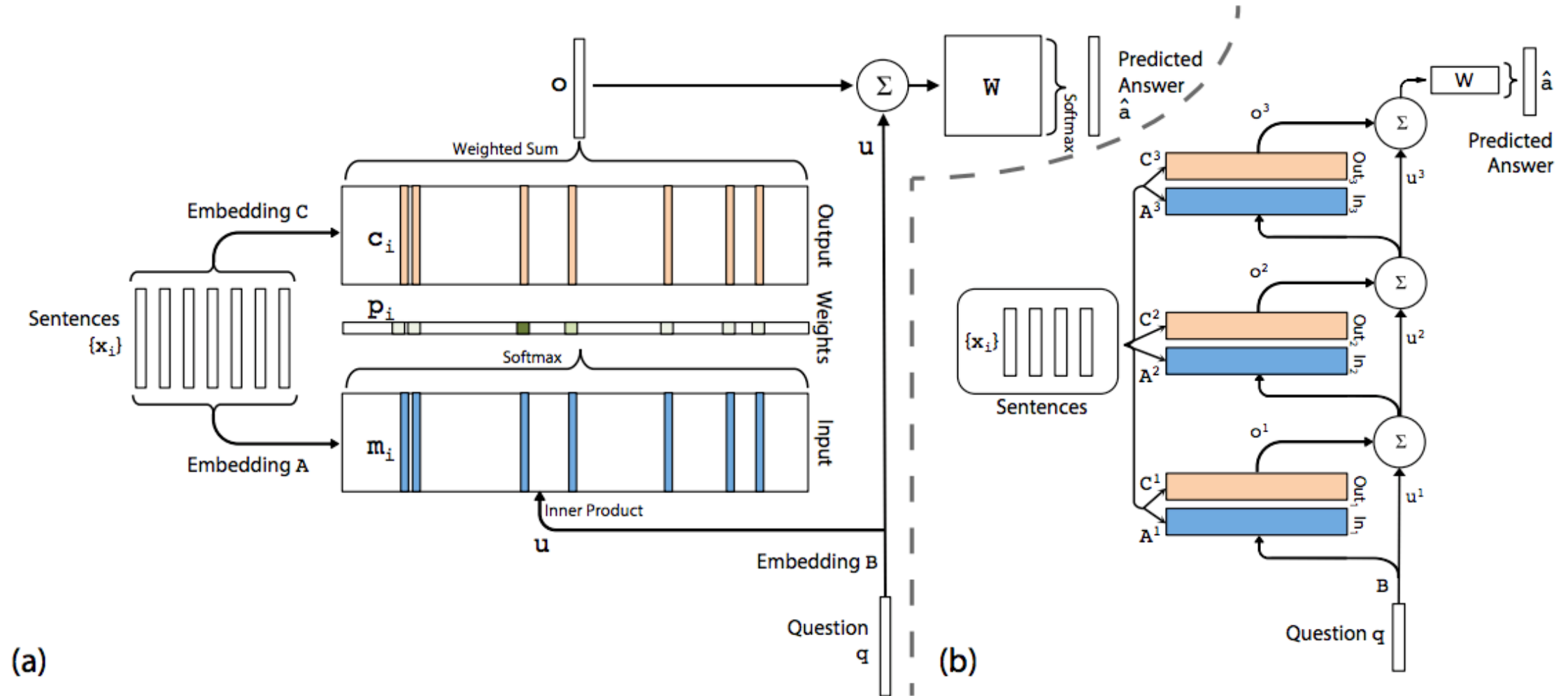
Impatient Attentive Reader

Considered transitive inference (bAbl)

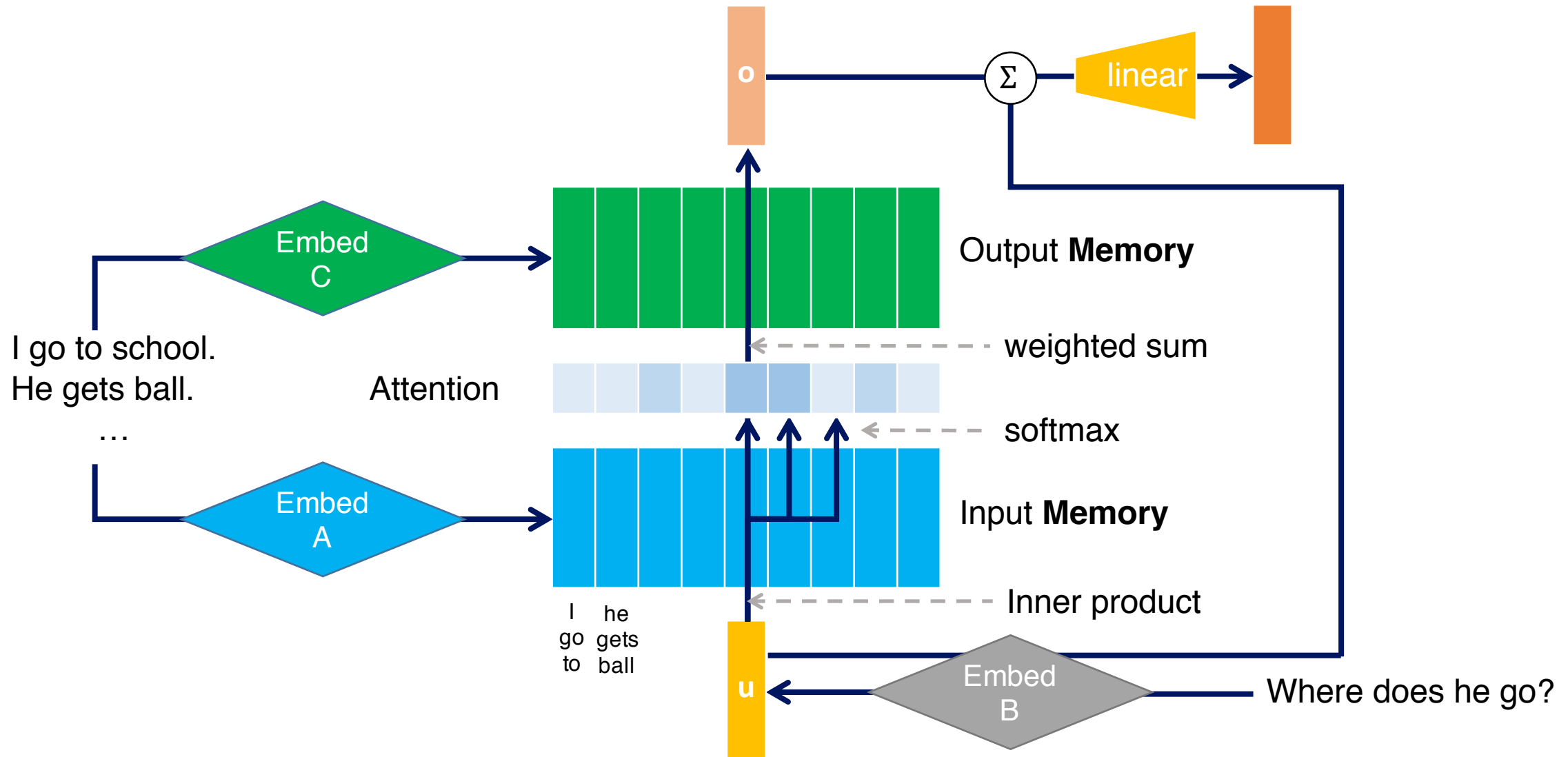
E2E MN

DMN

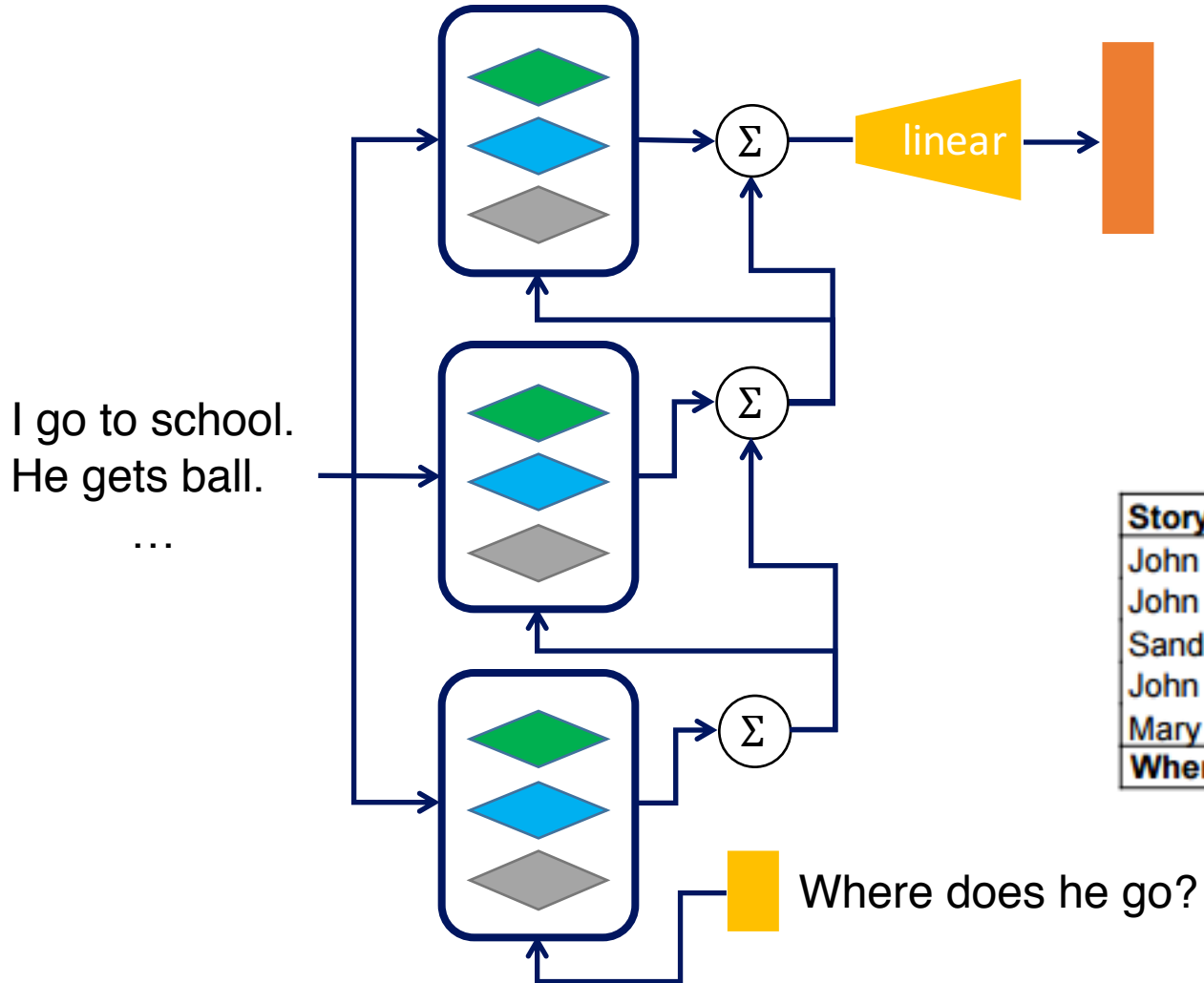
End-to-End Memory Network [Sukhbaatar, 2015]



End-to-End Memory Network [Sukhbaatar, 2015]



End-to-End Memory Network [Sukhbaatar, 2015]



Sentence representation :

i th sentence : $x_i = \{x_{i1}, x_{i2}, \dots, x_{in}\}$

BoW : $m_i = \sum_j A x_{ij}$

Position Encoding : $m_i = \sum_j l_j \cdot A x_{ij}$

Temporal Encoding : $m_i = \sum_j A x_{ij} + T_A(i)$

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				

Training details

Linear Start (LS) help avoid local minima

- First train with softmax in each memory layer removed, making the model entirely linear except for the final softmax
- When the validation loss stopped decreasing, the softmax layers were re-inserted and training recommenced

RNN-style layer-wise weight tying

- The input and output embeddings are the same across different layers

Learning **time invariance** by injecting random noise

- Jittering the time index with random empty memories
- Add “dummy” memories to regularize $T_A(i)$

Example of bAbl 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 (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 (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 (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				

The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations [Hill, 2016]

S: 1 Mr. Cropper was opposed to our hiring you .
2 Not , of course , that he had any personal objection to you , but he is set
against female teachers , and when a Cropper is set there is nothing on earth can
change him .
3 He says female teachers ca n't keep order .
4 He 's started in with a spite at you on general principles , and the boys know
it .
5 They know he 'll back them up in secret , no matter what they do , just to prove
his opinions .
6 Cropper is sly and slippery , and it is hard to corner him . ''
7 `` Are the boys big ? ''
8 queried Esther anxiously .
9 `` Yes .
10 Thirteen and fourteen and big for their age .
11 You ca n't whip 'em -- that is the trouble .
12 A man might , but they 'd twist you around their fingers .
13 You 'll have your hands full , I 'm afraid .
14 But maybe they 'll behave all right after all . ''
15 Mr. Baxter privately had no hope that they would , but Esther hoped for the
best.
16 She could not believe that Mr. Cropper would carry his prejudices into a
personal application .
17 This conviction was strengthened when he overtook her walking from school the
next day and drove her home .
18 He was a big , handsome man with a very suave , polite manner .
19 He asked interestedly about her school and her work , hoped she was getting on
well , and said he had two young rascals of his own to send soon .
20 Esther felt relieved .

Q: She thought that Mr. _____ had exaggerated matters a little .

C: Baxter, Cropper, Esther, course, fingers, manner, objection, opinion, right, spite.

a: Baxter

The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations [Hill, 2016]

- Context sentences : $S = \{s_1, s_2, \dots, s_n\}$, s_i : BoW word representation
- Encoded memory : $m_i = \phi(s) \forall s \in S$
- Lexical memory
 - Each word occupies a separate slot in the memory
 - s is a single word and $\phi(s)$ has only one non-zero feature
 - Multiple hop only beneficial in this memory model
- **Window memory (best)**
 - s corresponds to a window of text from the context S centered on an individual mention of a candidate c in S
$$m_i = \{w_{i-(b-1)/2} \dots w_i \dots w_{i+(b-1)/2}\}$$
 - Where $w_i \in C$ which is an instance of one of the candidate words
- Sentential memory
 - Same as original implementation of Memory Network

The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations [Hill, 2016]

Self-supervision for window memories

- Memory supervision (knowing which memories to attend to) is not provided at training time
- Making gradient steps using SGD to **force** the model to give a **higher score to the supporting memory \tilde{m}** relative to any other memory from any other candidate using:

$$\text{Hard attention (training and testing)} : m_{o1} = \underset{i=1,\dots,n}{\operatorname{argmax}} c_i^T q$$

$$\text{Soft attention (testing)} : m_{o1} = \sum_{i=1\dots n} \alpha_i m_i, \text{ with } \alpha_i = \frac{e^{c_i^T q}}{\sum_j e^{c_j^T q}}$$

- If m_{o1} happens to be different from \tilde{m} (memory contain true answer), then model is updated
- Can be understood as **a way of achieving *hard attention over memories*** (no need any new label information beyond the training data)

The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations [Hill, 2016]

S: 1 So they had to fall (a long way) .
2 So they got their tails fast (in their mouths) .
3 So they could n't get them out again .
4 That 's all . '
5 ` Thank you , ' said Alice , ` it 's very interesting .
6 I never knew so much (about a whiting before) . ' '
7 I can tell you more than that , if you like , ' said the Gryphon .
8 ` Do you know why it 's (called a whiting ?) ' '
9 I never thought about it , ' said Alice .
10 ` Why ? '
11 ` IT (DOES THE BOOTS AND SHOES) . '
12 the Gryphon replied very solemnly .
13 Alice was thoroughly puzzled .
14 ` (Does the boots and shoes) ! '
15 she repeated in a wondering tone .
16 ` Why , what (are YOUR shoes done with) ? '
17 said the Gryphon . '
18 I mean , what makes them so shiny ? '
19 Alice looked down at them , and considered a little before she (gave
her answer) .
20 They 're done with blacking , I believe .

Q: `Boots and shoes under the sea , ' the _____ went on in a deep voice , are done (with a whiting) .

C: Alice, BOOTS, Gryphon, SHOES, answer, fall, mouths, tone, way, whiting.

MemNNs (window + self-sup.): **Gryphon**

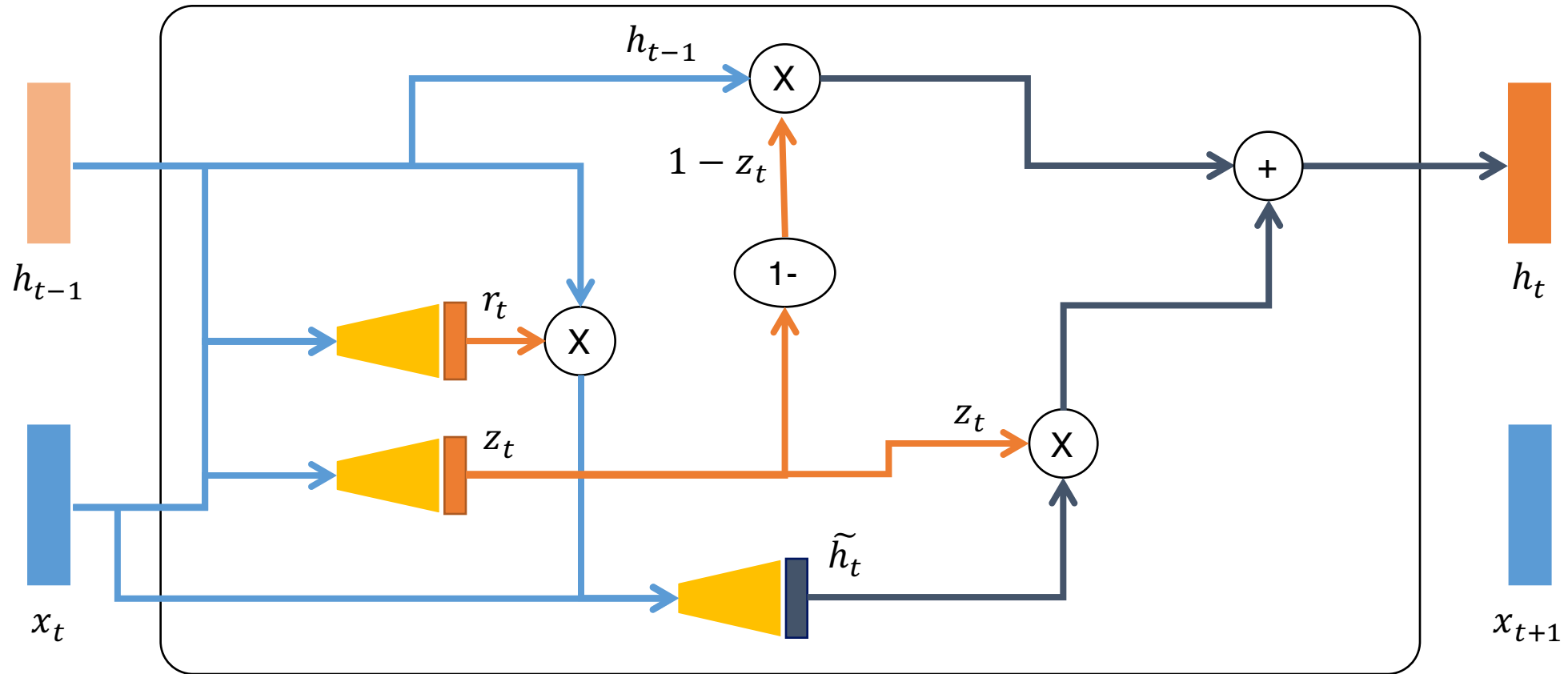
S: 1 He thought that Old Mr. Toad was trying to fool him .
2 Presently Peter Rabbit came along .
3 He found Jimmy Skunk sitting in a brown study .
4 He had quite forgotten to look for fat beetles , and when he (forgets to do
that you) may make up your mind that Jimmy is doing some hard thinking .
5 `` Hello , old Striped-coat , what have you got on your mind this fine
morning ? ''
6 cried Peter Rabbit .
7 `` Him , '' said Jimmy simply , pointing down the Lone Little Path .
8 Peter looked .
9 `` (Do you mean) Old Mr. Toad ! ''
10 he asked .
11 Jimmy nodded .
12 `` (Do you see) anything queer about him ? ''
13 he asked in his turn .
14 `` (Do you see) anything queer about him ? ''
15 he asked .
16 Peter stared down the Lone Little Path .
17 `` No , '' he replied , `` except that he seems in a great hurry . ''
18 `` That 's just it , '' Jimmy returned promptly .
19 `` Did (you ever see him hurry) unless he was frightened ? ''
20 Peter confessed that he never had

Q: `` Well , he is n't _____ now , yet just look at him go '' retorted Jimmy .

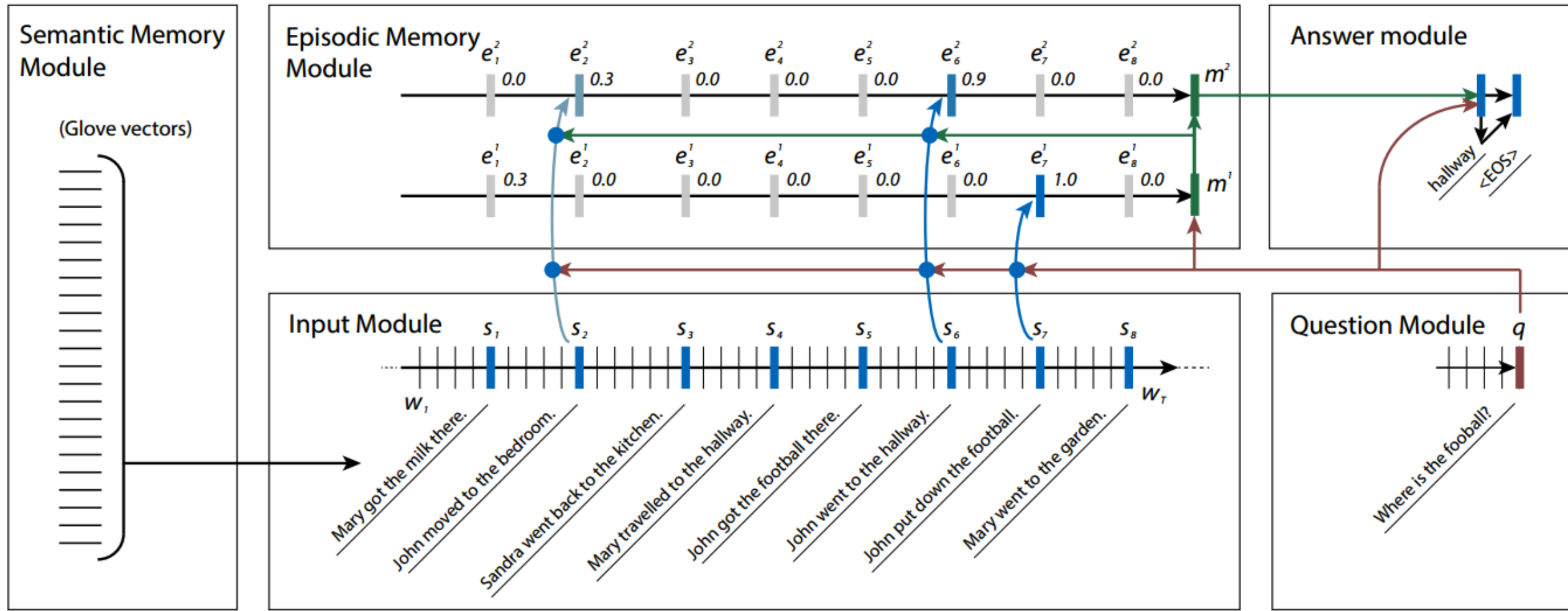
C: Do, came, confessed, frightened, mean, replied, returned, said, see, thought.

MemNNs (window +self-sup.): **frightened**

Gated Recurrent Network (GRU)

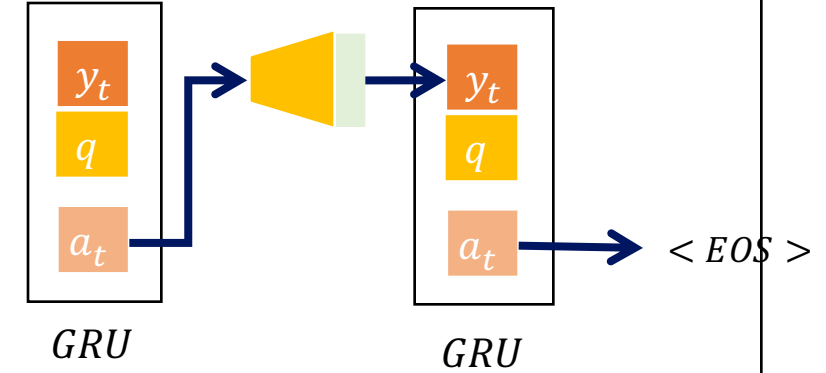
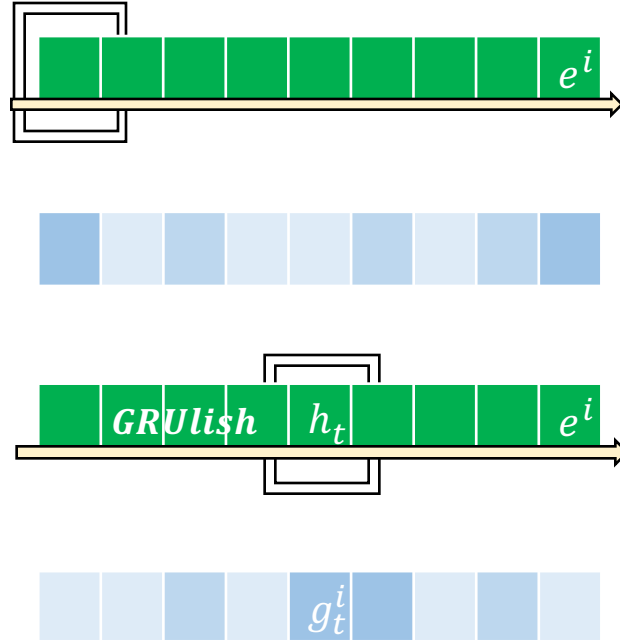


Ask Me Anything: Dynamic Memory Networks for Natural Language Processing [Kumar, 2015]



Ask Me Anything: Dynamic Memory Networks for Natural Language Processing [Kumar, 2015]

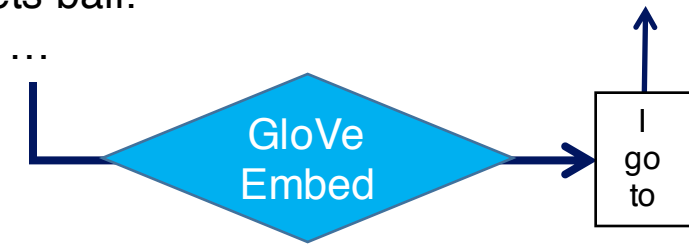
Episodic Memory



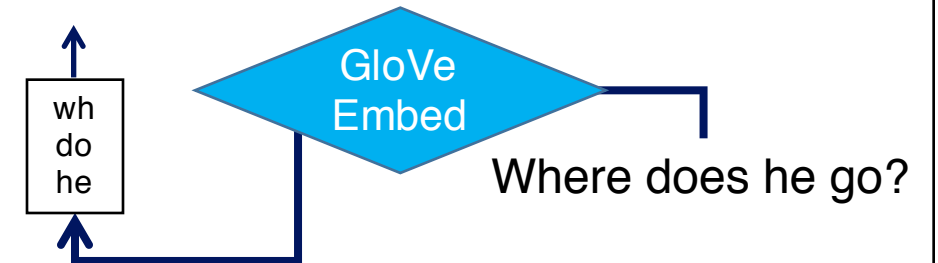
Answer Module

I go to school.
He gets ball.
...

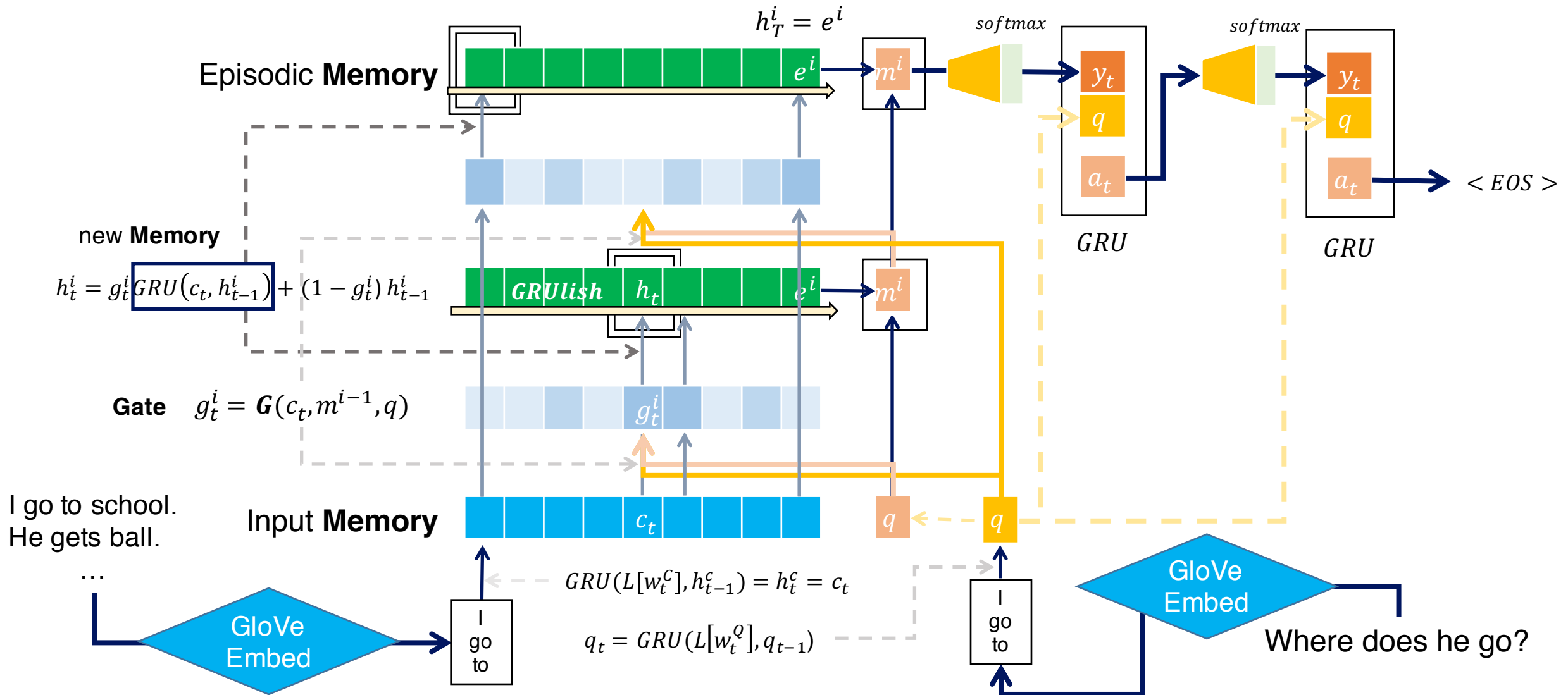
Input Module



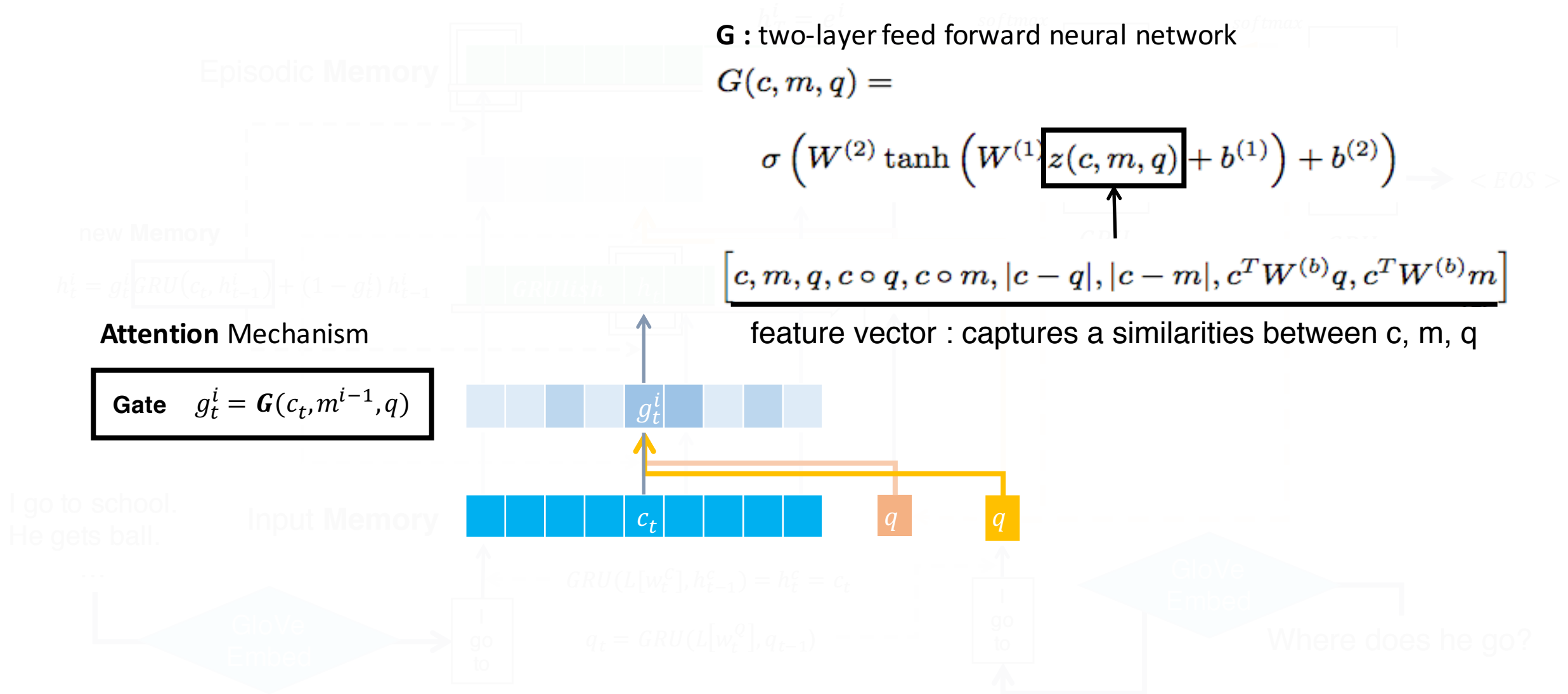
Question Module



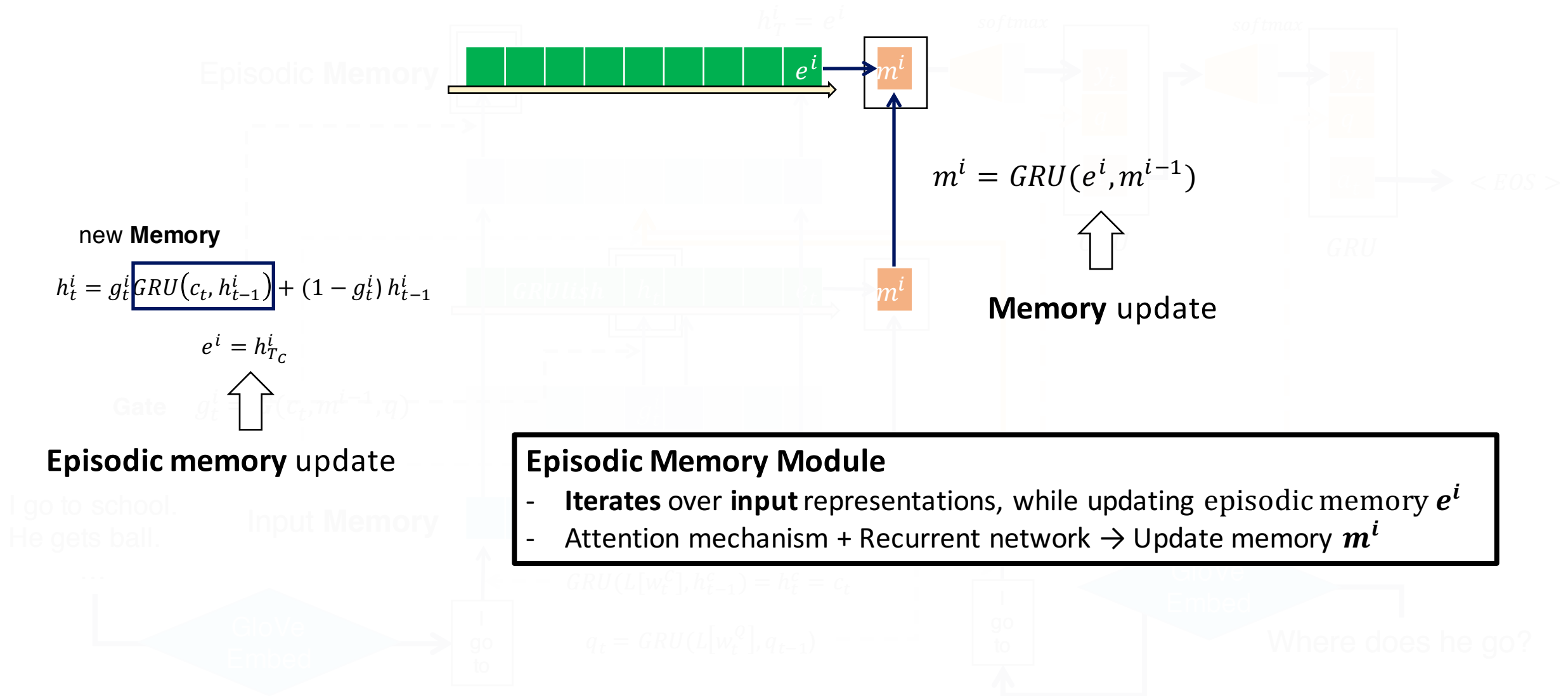
Ask Me Anything: Dynamic Memory Networks for Natural Language Processing [Kumar, 2015]



Ask Me Anything: Dynamic Memory Networks for Natural Language Processing [Kumar, 2015]



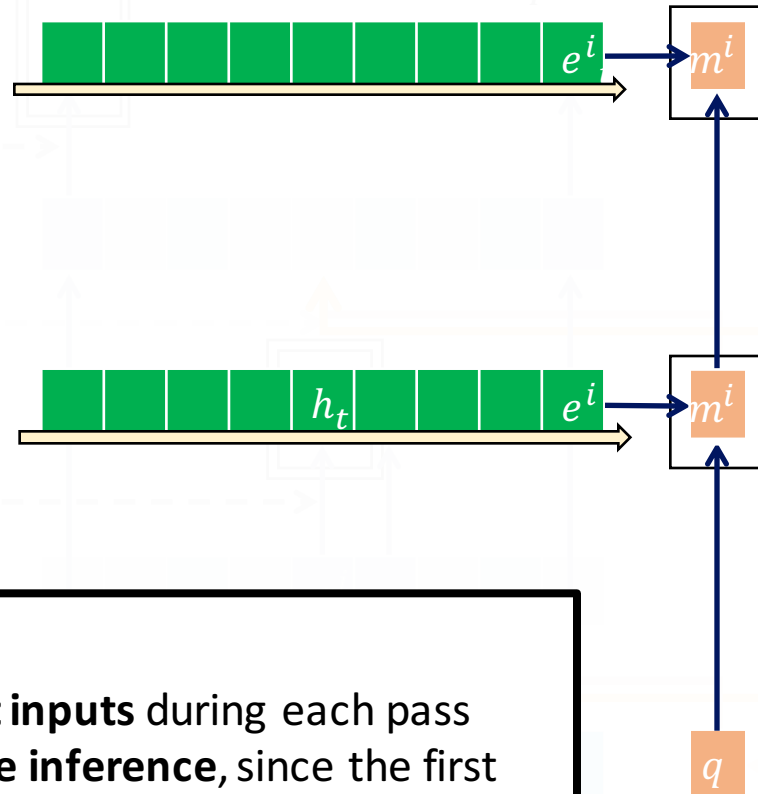
Ask Me Anything: Dynamic Memory Networks for Natural Language Processing [Kumar, 2015]



Ask Me Anything: Dynamic Memory Networks for Natural Language Processing [Kumar, 2015]

Max passes	task 3 three-facts	task 7 count	task 8 lists/sets	sentiment (fine grain)
0 pass	0	48.8	33.6	50.0
1 pass	0	48.8	54.0	51.5
2 pass	16.7	49.1	55.6	52.1
3 pass	64.7	83.4	83.4	50.1
5 pass	95.2	96.9	96.5	N/A

Table 4. Effectiveness of episodic memory module across tasks. Each row shows the final accuracy in term of percentages with a different maximum limit for the number of passes the episodic memory module can take. Note that for the 0-pass DMN, the network essentially reduces to the output of the attention module.



Criteria for Stopping

- Append a special end-of-passes representation to the input c
- Stop if this representation is **chosen** by the **gate** function
- Set a maximum number of iterations
- This is why called **Dynamic** MM

Multiple Episodes

- Allows to **attend** to **different inputs** during each pass
- Allows for a type of **transitive inference**, since the first pass may uncover the need to retrieve additional facts.

Q : Where is the football?

C1 : John put down the football.

Only once the model sees C1, John is relevant, can reason that the second iteration should retrieve where John was.

Training Details

- Adam optimization
- L_2 regularization, dropout on the word embedding (GloVe)

bAbI dataset

- Objective function : $J = \alpha E_{CE}(Gates) + \beta E_{CE}(Answers)$
- **Gate supervision** aims to select **one sentence per pass**
 - Without supervision : GRU of c_t, h_t^i and $e^i = h_{T_C}^i$
 - With supervision (simpler) : $e^i = \sum_{t=1}^T softmax(g_t^i) c_t$, where $softmax(g_t^i) = \frac{\exp(g_t^i)}{\sum_{j=1}^T \exp(g_j^i)}$ and g_t^i is the value before sigmoid
 - Better results, because softmax encourages **sparsity** & suited to **picking one** sentence

Training Details

Stanford Sentiment Treebank (Sentiment Analysis)

- Use all full sentences, subsample 50% of phrase-level labels every epoch
- Only evaluated on the full sentences
- Binary classification, neutral phrases are removed from the dataset
- Trained with GRU sequence models

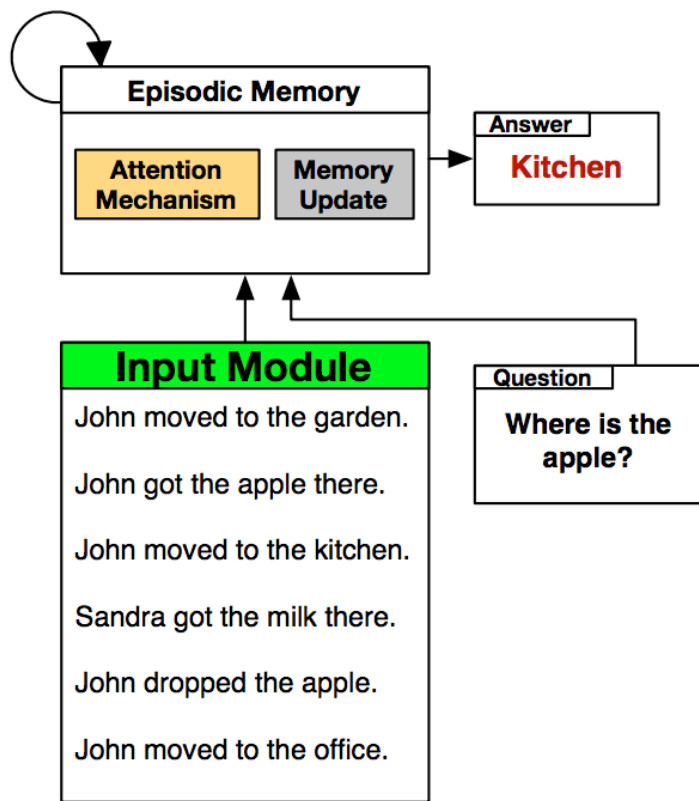
Task	Binary	Fine-grained
MV-RNN	82.9	44.4
RNTN	85.4	45.7
DCNN	86.8	48.5
PVec	87.8	48.7
CNN-MC	88.1	47.4
DRNN	86.6	49.8
CT-LSTM	88.0	51.0
DMN	88.6	52.1

Training Details

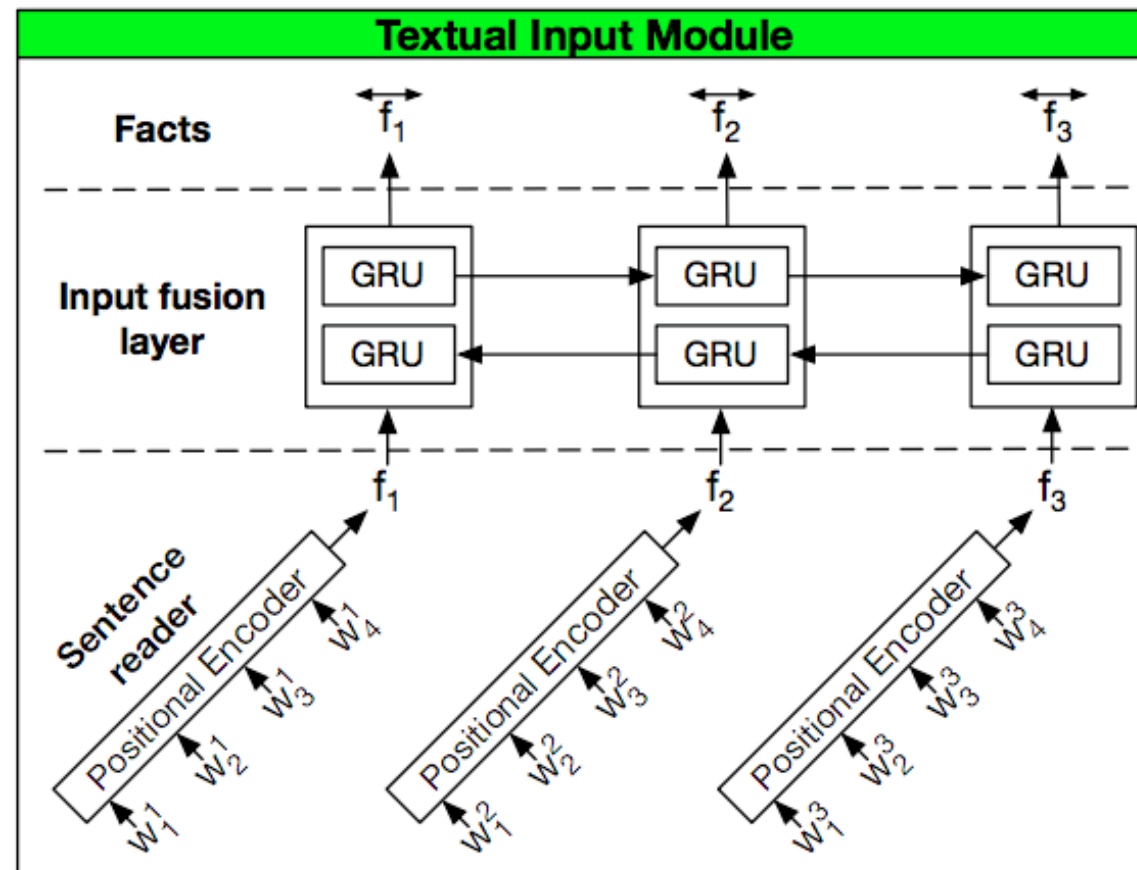
Question: Where was Mary before the Bedroom?
Answer: Cinema.

Facts	Episode 1	Episode 2	Episode 3
Yesterday Julie traveled to the school.			
Yesterday Marie went to the cinema.			
This morning Julie traveled to the kitchen.			
Bill went back to the cinema yesterday.			
Mary went to the bedroom this morning.			
Julie went back to the bedroom this afternoon.			
[done reading]			

Dynamic Memory Networks for Visual and Textual Question Answering [Xiong 2016]



(a) Text Question-Answering



Several design choices are **motivated by intuition** and **accuracy improvements**

Input Module for DMN

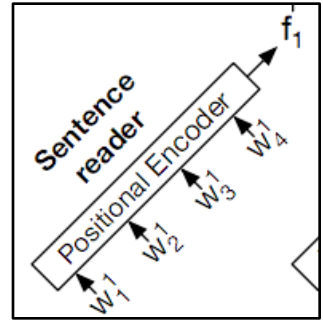
- A single GRU for embedding story and store the hidden states
- GRU provides **temporal component** by allowing a sentence to know the **content of** the sentences that came **before them**
- **Cons:**
 - GRU only allows sentences to have context from sentences **before** them, but **not after them**
 - **Supporting sentences** may be too **far** away from each other
- Here comes **Input fusion** layer

Input Module for DMN+

Replacing a single GRU with two different components

1. Sentence reader : responsible only for encoding the **words into a sentence embedding**

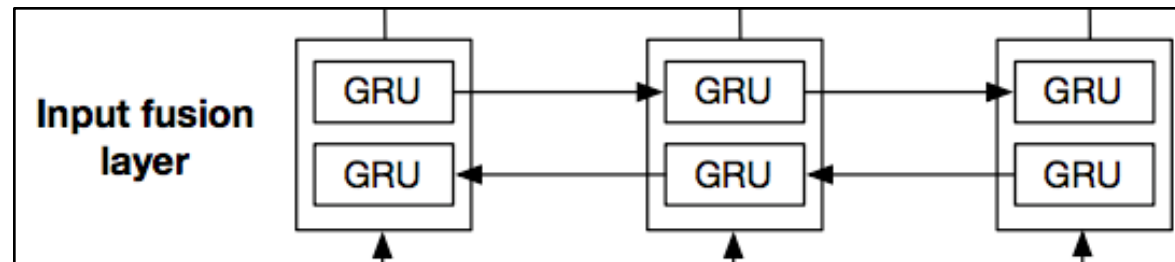
- Use positional encoder (used in E2E) : $f_i = \sum_j l_j \cdot Ax_{ij}$
- Considered GRUs LSTMs, but required more computational resources, prone to overfitting



2. Input fusion layer : interactions between sentences, allows **content interaction** between sentences

- **bi-directional** GRU to allow information from both past and future sentences
- gradients do not need to propagate through the words between sentences
- **distant supporting sentences** can have a more **direct interaction**

$$\begin{aligned}\vec{f}_i &= GRU_{fwd}(f_i, \vec{f}_{i-1}) \\ \overleftarrow{f}_i &= GRU_{bwd}(f_i, \overleftarrow{f}_{i+1}) \\ \overleftrightarrow{f}_i &= \overleftarrow{f}_i + \vec{f}_i\end{aligned}$$



Input Module for DMN+

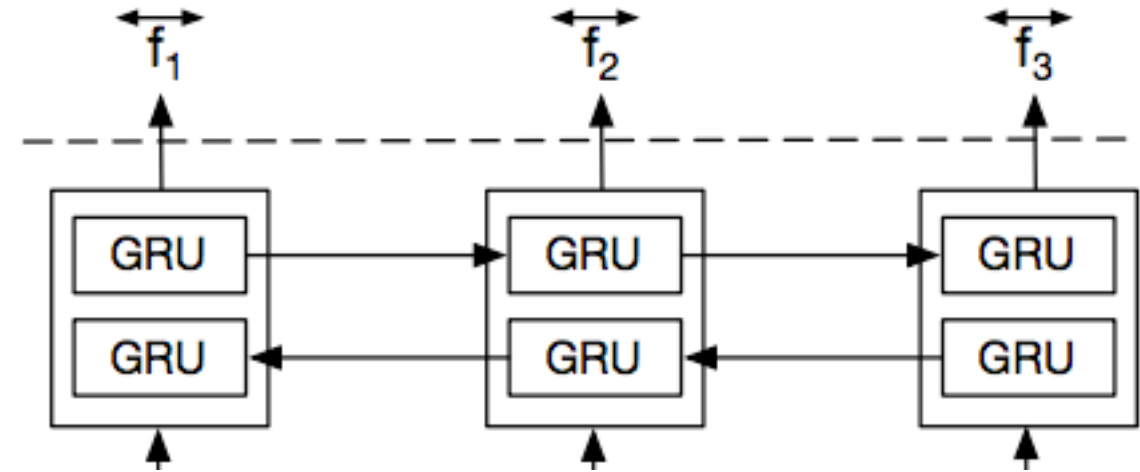
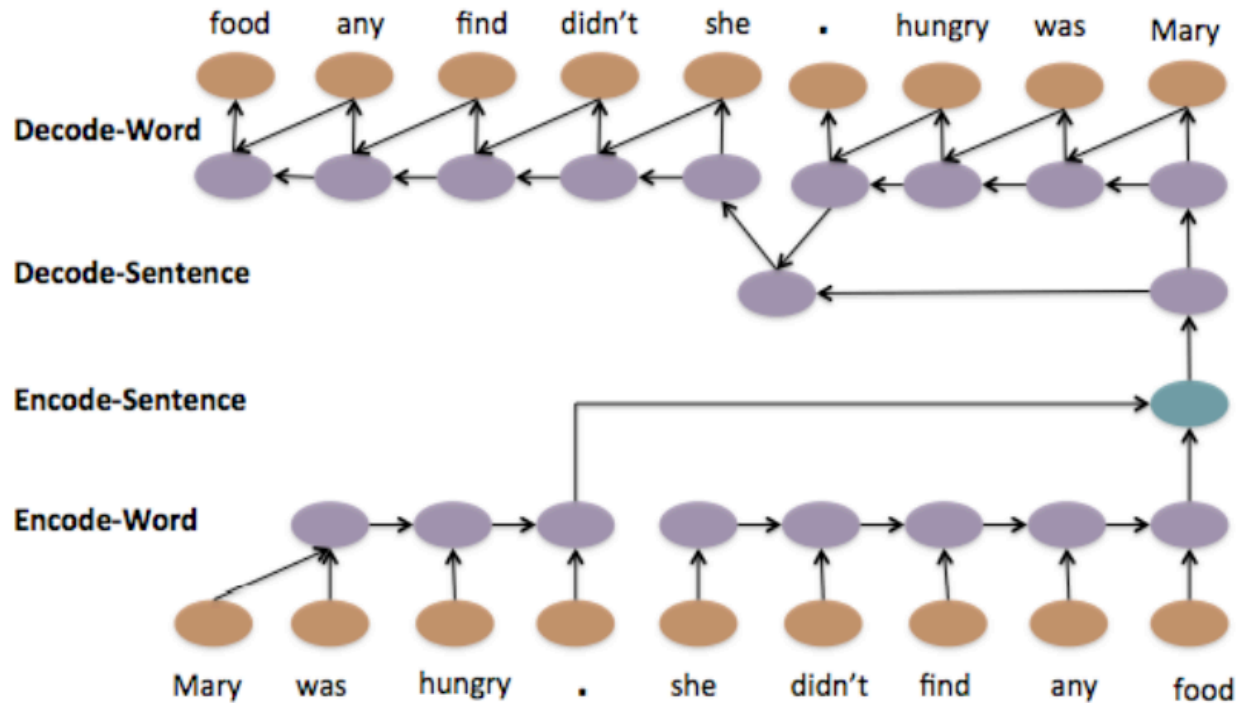


Figure 2: Hierarchical Sequence to Sequence Model.

Referenced paper : [A Hierarchical Neural Autoencoder for Paragraphs and Documents](#) [Li, 2015]

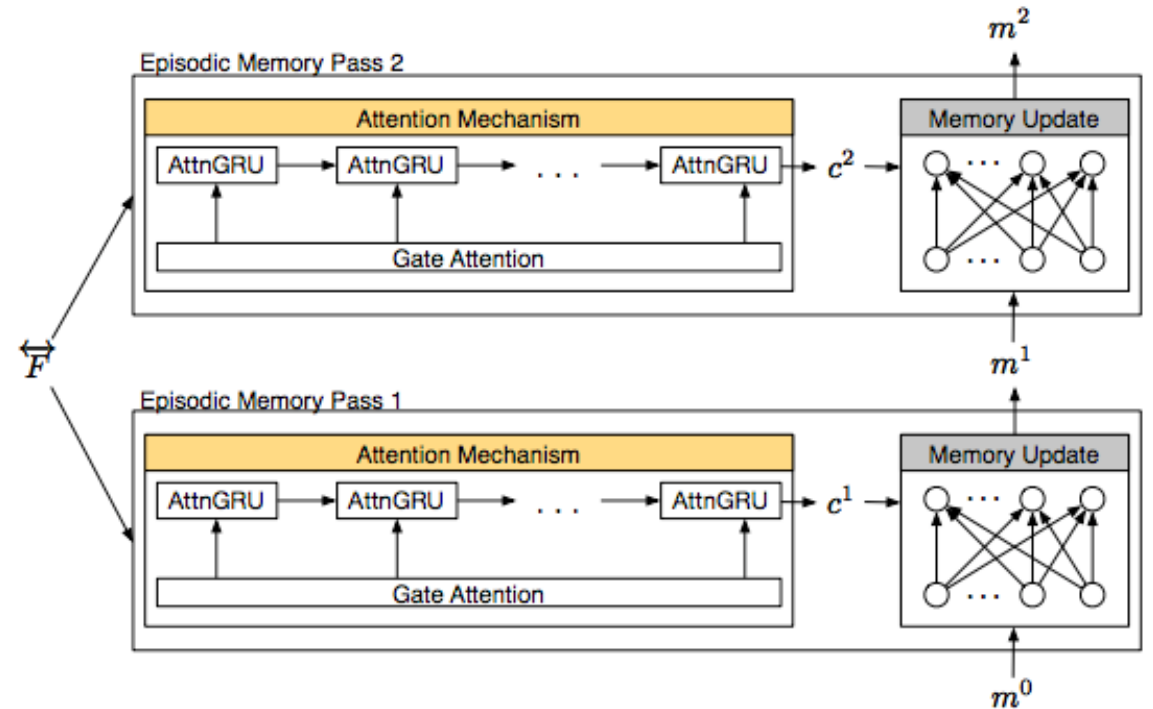
Episodic Memory Module for DMN+

- $\overleftrightarrow{F} = [\overleftrightarrow{f}_1, \overleftrightarrow{f}_2, \dots, \overleftrightarrow{f}_N]$: output of the input module
- interactions between the fact \overleftrightarrow{f}_i and both the question q and episode memory state m^t

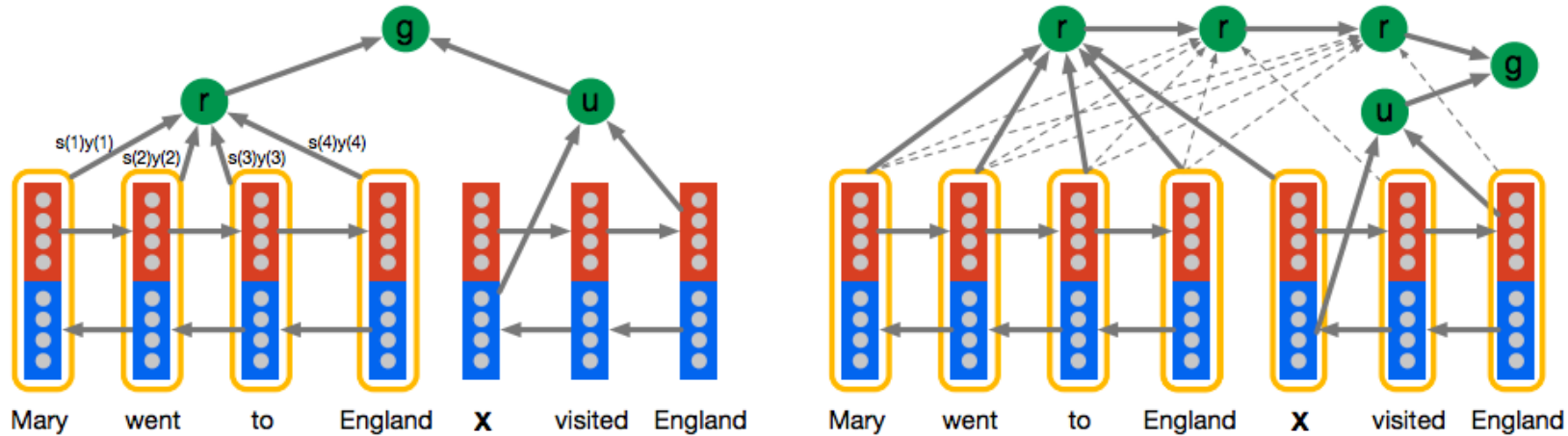
$$z_i^t = [\overleftrightarrow{f}_i \circ q; \overleftrightarrow{f}_i \circ m^{t-1}; |\overleftrightarrow{f}_i - q|; |\overleftrightarrow{f}_i - m^{t-1}|]$$

$$Z_i^t = W^{(2)} \tanh \left(W^{(1)} z_i^t + b^{(1)} \right) + b^{(2)}$$

$$g_i^t = \frac{\exp(Z_i^t)}{\sum_{k=1}^{M_i} \exp(Z_k^t)}$$

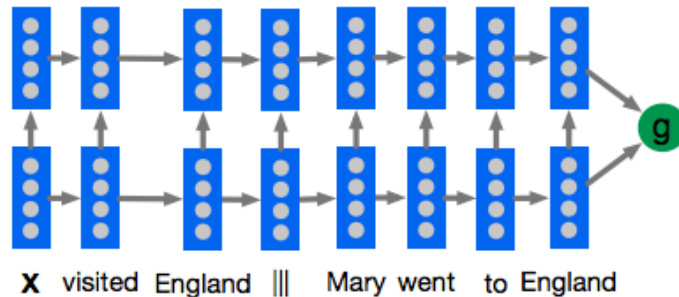


Teaching Machines to Read and Comprehend [Hermann 2015]



(a) Attentive Reader.

(b) Impatient Reader.



(c) A two layer Deep LSTM Reader with the question encoded before the document.

Neural Variational Inference for Text Processing [Miao, 2015]

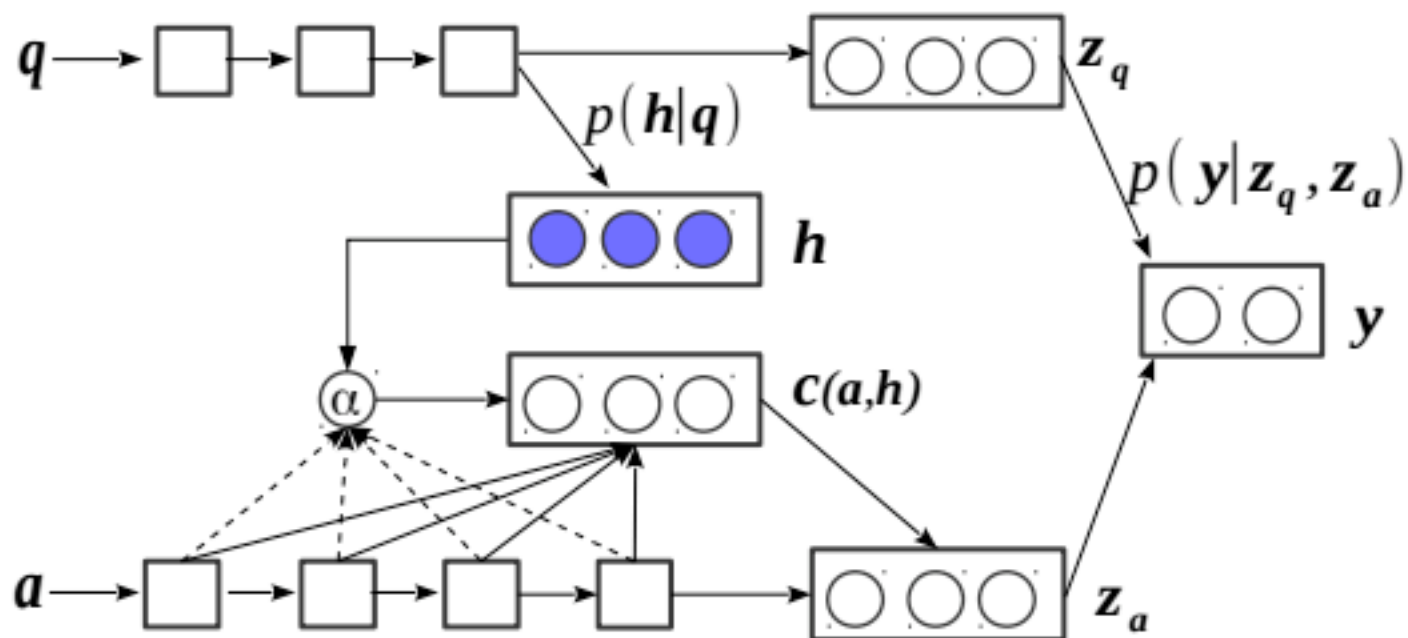
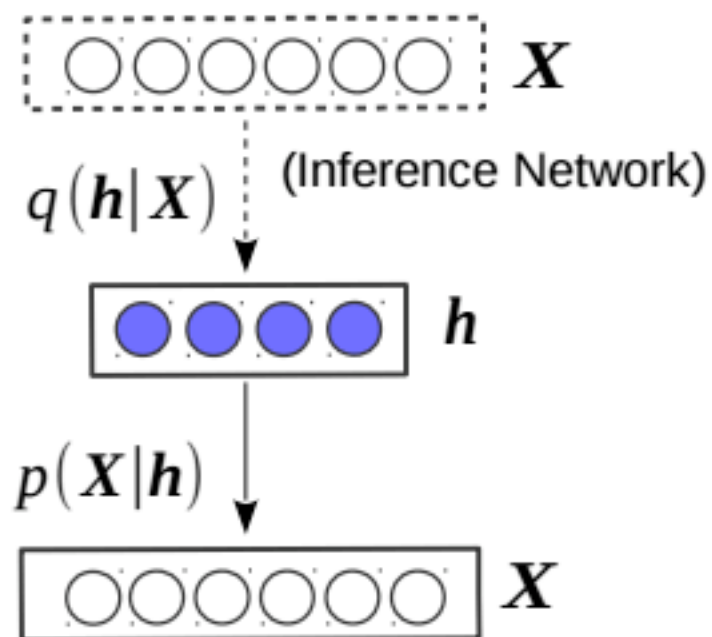


Figure 1: NVDM for document modelling.

Figure 2: NASM for question answer selection.

Stochastic Latent Variable

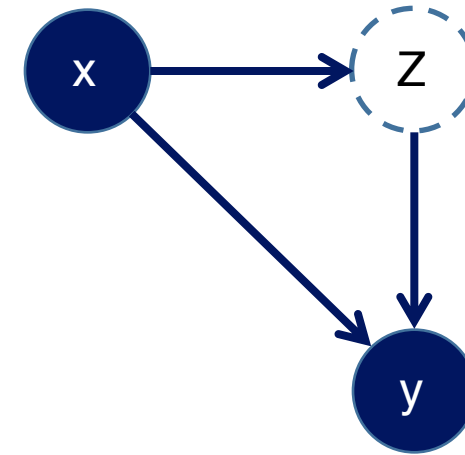
Generative Model



$$p(x) = \sum_z p(x, z) = \sum_z p(x|z)p(z)$$

$$p(x) = \int_z p(x, z) = \int_z p(x|z)p(z)$$

Conditional Generative Model



$$p(y|x) = \sum_z p(y|z, x)p(z|x)$$

$$p(y|x) = \int_z p(y|z, x)p(z|x)$$

Variational Inference Framework

$$p(x, z) = p(x|z)p(z) = \sum_h p(x|h)p(h|z)p(z)$$

$$\log p_\theta(x, z) = \log \int_h \frac{q(h)}{q(h)} p(x|h)p(h|z)p(z) dh \geq \int_h q(h) \log \frac{p(x|h)p(h|z)p(z)}{q(h)} dh$$

$$= \int_h q(h) \log \frac{p(x|h)p(h|z)}{q(h)} dh + \int_h q(h) \log \frac{p(z)}{q(h)} dh$$

$$= E_{q(h)}[\log p(x|h)p(h|z) - \log q(h)] - D_{KL}(q(h)||p(z))$$

$$= E_{q(h)}[\log p(x|h)p(h|z)p(z) - \log q(h)]$$

Variational Inference Framework

$$p_{\theta}(x, z) = p_{\theta}(x|z)p(z) = \sum_h p_{\theta}(x|h)p_{\theta}(h|z)p(z)$$

Jensen's Inequality

$$\log p_{\theta}(x, z) = \log \int_h \frac{q(h)}{q(h)} p_{\theta}(x|h)p_{\theta}(h|z)p(z)dh \geq \int_h q(h) \log \frac{p_{\theta}(x|h)p_{\theta}(h|z)p(z)}{q(h)} dh$$

$$= \int_h q(h) \log \frac{p_{\theta}(x|h)p_{\theta}(h|z)}{q(h)} dh + \int_h q(h) \log \frac{p(z)}{q(h)} dh$$

$$= E_{q(h)}[\log p_{\theta}(x|h)p_{\theta}(h|z) - \log q(h)] - D_{KL}(q(h)||p(z))$$

$$= E_{q(h)}[\log p_{\theta}(x|h)p_{\theta}(h|z) - \log q(h)] \quad \text{a tight lower bound if } q(h) = p(h|x, z)$$

Conditional Variational Inference Framework

$$p_{\theta}(y|x) = \sum_z p_{\theta}(y, z|x) = \sum_z p_{\theta}(y|x, z)p_{\pi}(z|x)$$

Jensen's Inequality

$$\log p(y|x) = \log \int_z \frac{q(z)}{q(z)} p(y|z, x) p(z|x) dz \geq \int_z q(z) \log \frac{p(y|z, x) p(z|x)}{q(z)} dz$$

$$= \int_z q(z) \log \frac{p(y|z, x)}{q(z)} dz + \int_h q(z) \log \frac{p(z|x)}{q(z)} dz$$

$$= \int_z q(z) \log p(y|z, x) dz - \int_z q(z) \log q(z) dz + \int_h q(z) \log \frac{p(z|x)}{q(z)} dz$$

$$= E_{q(z)}[\log p(y|z, x) - \log q(z)] - D_{KL}(q(z) \parallel p(z|x))$$

$$= E_{q(z)}[\log p(y|z, x) - \log q(z)] \quad \text{a tight lower bound if } q(z) = p(z|x)$$

Neural Variational Inference Framework

$$\log p_{\theta}(x, z) \geq E_{q(z)}[\log p(y|z, x) - \log q(z)] - D_{KL}(q(z) \parallel p(z|x)) = \mathcal{L}$$

1. Vector representations of the observed variables

$$u = f_z(z), v = f_x(x)$$

2. Joint representation (concatenation)

$$\pi = g(u, v)$$

3. Parameterize the variational distribution

$$\mu = l_1(\pi), \sigma = l_2(\pi)$$

Neural Variational Document Model

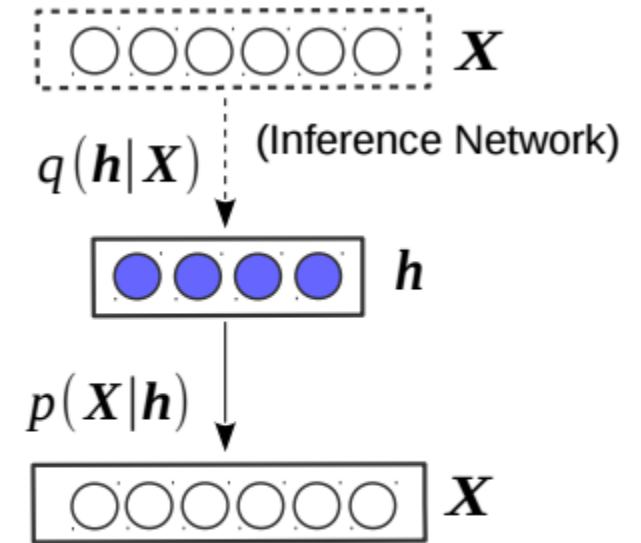


Figure 1: NVDM for document modelling.