# Memory Networks



Jetic Gū

# Overview

- Focus: Neural Knowledge Incorporation (in NLP)

- Architecture: Dedicated Memory Component in Neural Model

- Core Ideas:

  1. Background: Seq2Seq, Attention, etc.

  2. Memory Network

  3. Applications of Memory Network

  4. Future Work

# Review, and Limitations of Seq2Seq

Including Transformer, BERT, etc.

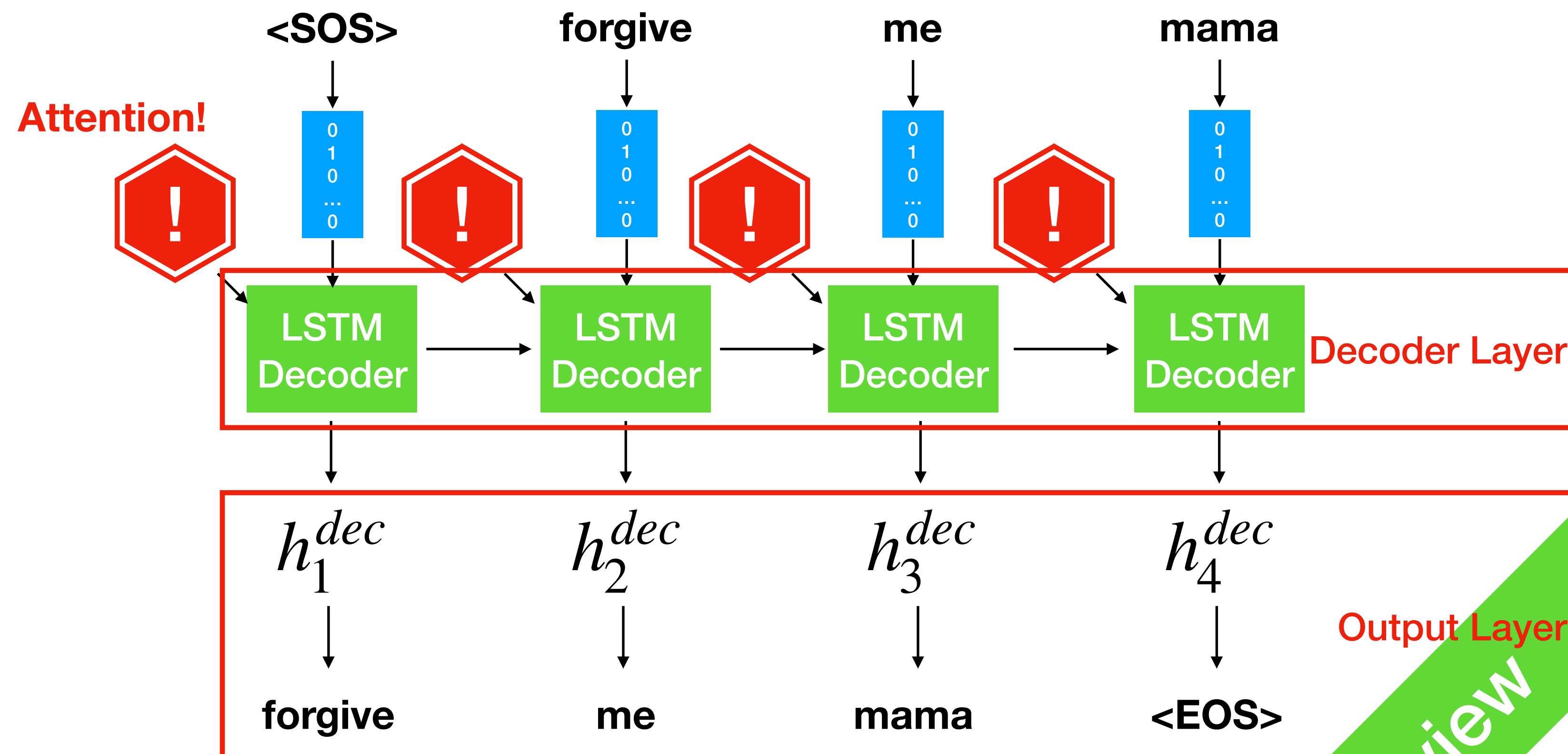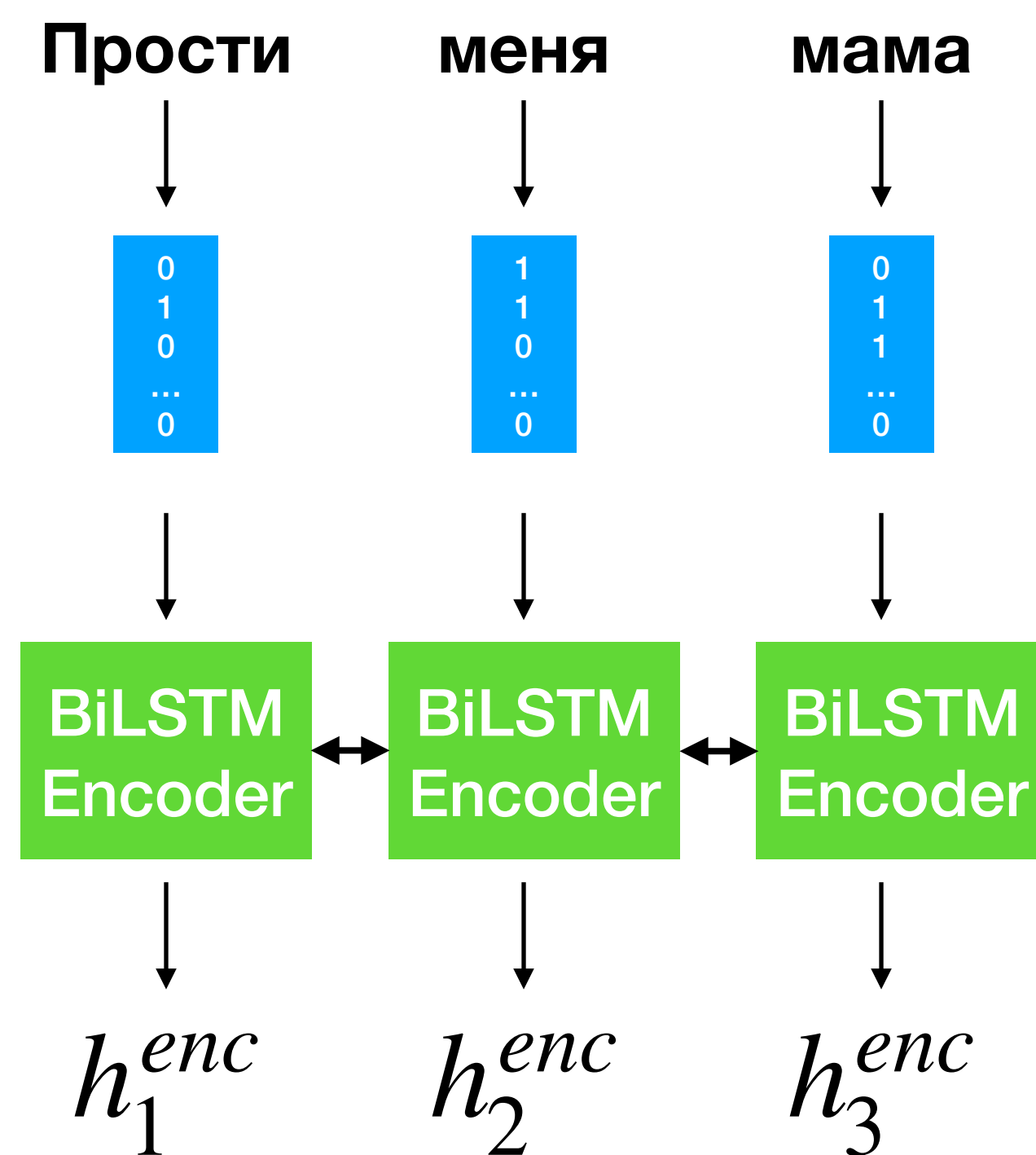**P1
Why?**

# Think: Seq2Seq Models

- Powerful

  - Variable-length input: variable-length output

  - Variants: Attentional Seq2Seq, Transformer, etc.

- Uses encoder decoder architecture

  - Encoder: performs feature extraction

  - Decoder: generates output based on aggregated features

- Roll of neural network units: projection functions, from one feature space to another

**Review**

# Think: Seq2Seq Models

- As an example, LSTM-based Seq2Seq for Translation
  Русский (RU) to English (EN)

**Прости**   **меня**   **мама**

**<SOS>**   **forgive**   **me**   **mama**

**Attention!**

LSTM
Decoder   LSTM
Decoder   LSTM
Decoder   LSTM
Decoder   **Decoder Layer**

BiLSTM
Encoder ↔ BiLSTM
Encoder ↔ BiLSTM
Encoder

$h_1^{enc}$   $h_2^{enc}$   $h_3^{enc}$

$h_1^{dec}$   $h_2^{dec}$   $h_3^{dec}$   $h_4^{dec}$   **Output Layer**

**forgive**   **me**   **mama**   **<EOS>**

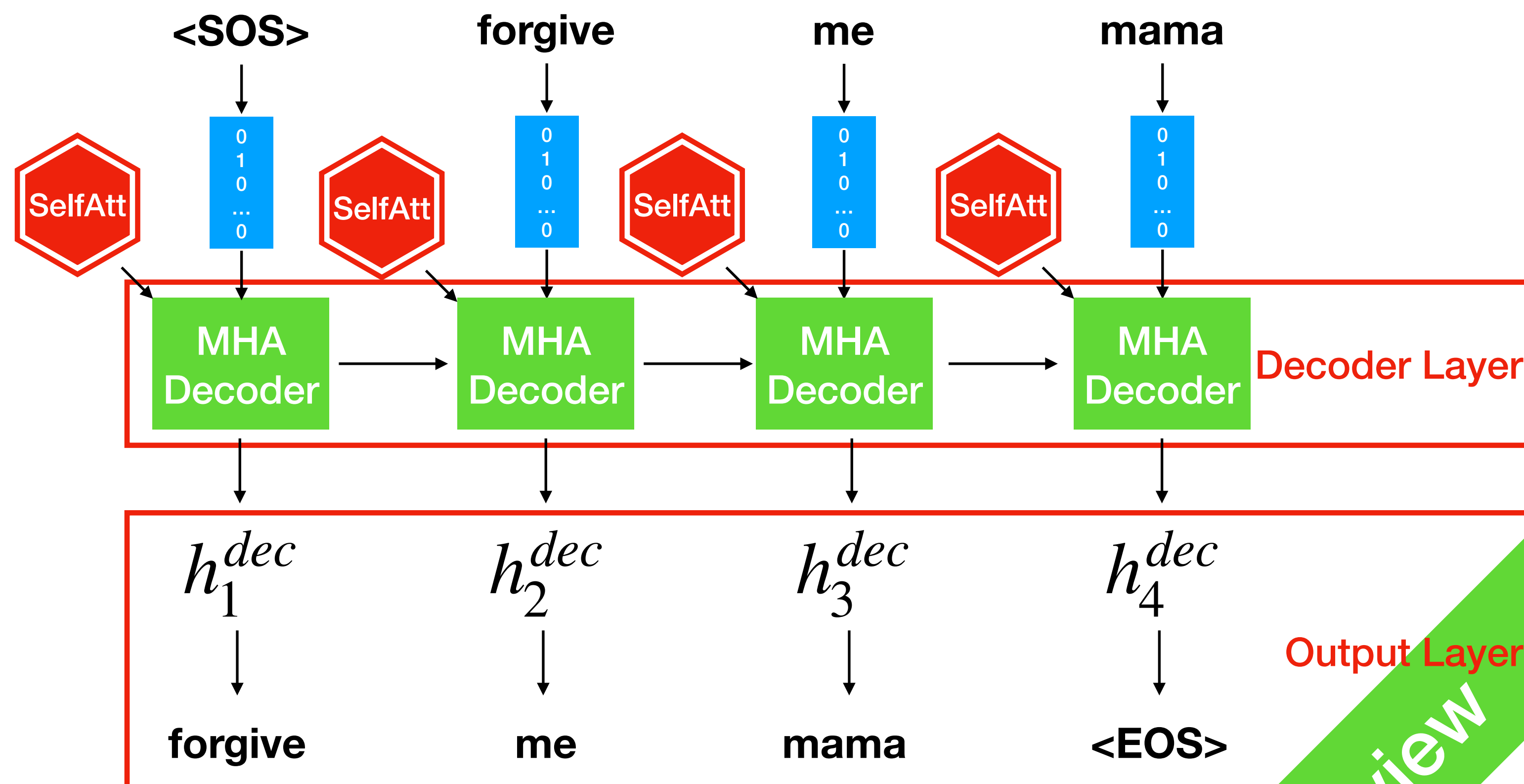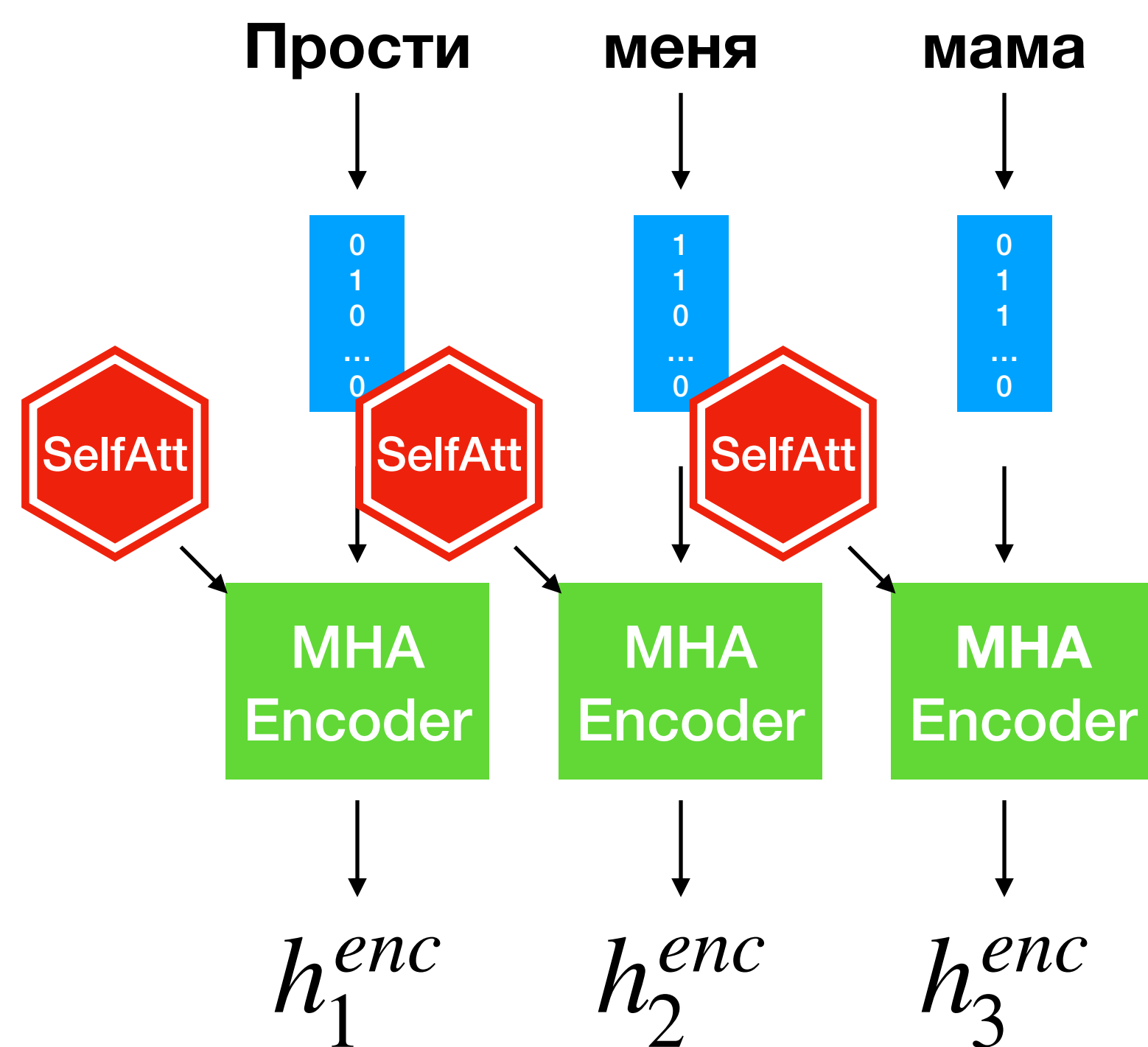**Review**

# Think: Transformer Models

**P1 Why?**

- As an example, Transformer-based Seq2Seq for Translation
  Русский (RU) to English (EN)

# Why does attention work?

- The neural decoder is a generative language model (think GPT)

  - $P(e_i \,|\, e_{<i}, F) \in [0,1]^{DictSize}$

  - Take $h^{enc}$ and $h^{dec}_{<i}$: provide condition for generation

    - $P(e_i \,|\, h^{dec}_{<i}, h^{enc}) = P(e_i \,|\, e_{<i}, f) \in [0,1]^{DictSize}$

    - What information is stored in $h^{enc}$ and $h^{dec}$
      And how do they contribute to the prediction of $e_i$?

**Think**

**P1 Why?** # Why does attention work?

- The neural decoder is a conditional language model

  - $P(e_i | e_{<i}, F) \in [0,1]^{DictSize}$

  - Guide with Attention for every step $t$

    - Project decoder state $h_t^{dec}$ and all encoder states $h_{0:|f|}^{enc}$ into the same feature space, find the most relevant $h_i^{enc}$ and do weighted sum

    - Decoder memory $|h^{dec}|$: ensure fluent language generation
      Encoder memory $|h^{enc}| \times |f|$: ensure src-tgt relevance

**Think**

**P1
Why?**

# What kind of Knowledge is learned in NMT for RU-EN?

- Russian Embedding and English Embedding: word-level features

  **Feature Map**

- Encoder: extract useful features from Russian Word Embeddings, w.r.t. global context

  **Features Aggregation & Representation Storage**

- Decoder: predict next word given previous English words & Encoder

  **Projection**

- Attention: project enc. rep. and previous dec. rep., aggregate the most relevant representations for the current time step

  **Projection**

- Output layer: project dec. rep. into target dictionary probability distribution

  **Projecti**

**Concept**

**P1
Why?**

# Current NLP Approaches

- **Symbolic knowledge** are fed into Neurones (e.g. RNN) for training

- **Limited** in memory **storage capacity**

- **Agnostic to explicit knowledge**, we assume the parameters will pick it up

- Applications: seq2seq, classification, etc.

Review

1. Graves et al., Hybrid computing using a neural network with dynamic external memory, Nature 2016

# RNN/MHA Units

Input *query* $\quad$ $x_0$ $\quad$ $x_1$ $\quad$ ... $\quad$ $x_{n-1}$

Review

# RNN/MHA Units

Input *query* | $x_0$ | $x_1$ | ... | $x_{n-1}$

| Embedding | Embedding | Embedding | Embedding |

| RNN/MHA | ↔ | RNN/MHA | ⇠ ⇢ | RNN/MHA |

Review

# RNN/MHA Units

Input *query*  $x_0$  $x_1$  ...  $x_{n-1}$

Embedding   Embedding   Embedding   Embedding

RNN/MHA  RNN/MHA  RNN/MHA

Neural Feature Representations  $h_0$  $h_1$  ...  $h_{n-1}$

Review

**P1
Why?**

# RNN/MHA Units

Input *query*

$x_0$   $x_1$   ...   $x_{n-1}$

Embedding   Embedding   Embedding   Embedding

RNN/MHA   RNN/MHA   RNN/MHA

Neural Feature Representations   $h_0$   $h_1$   ...   $h_{n-1}$

Neural Decoder

Review

# RNN/MHA Units

Input *query*

$x_0$    $x_1$    ...    $x_{n-1}$

Embedding    Embedding    Embedding    Embedding

RNN/MHA    RNN/MHA    RNN/MHA

Neural Feature Representations    $h_0$    $h_1$    ...    $h_{n-1}$

Neural Decoder

Response

**P1
Why?**
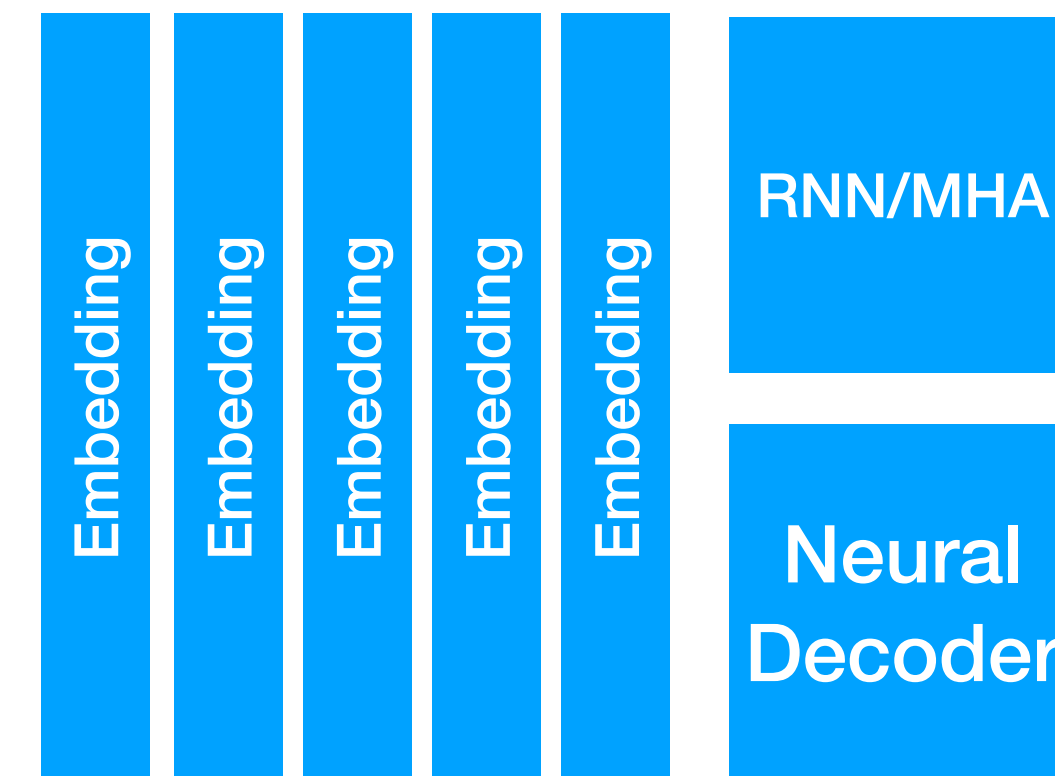
# Current NLP Approaches

- Neural Components

  - *Query*-to-*Response* Mapping Function

- Expect: **limited parameters** learn all knowledge

- Reality: sometimes you need information <u>external to the input</u>: ***Context***

Embedding Embedding Embedding Embedding Embedding RNN/MHA

**Neural Decoder**

Review

**P1
Why?**

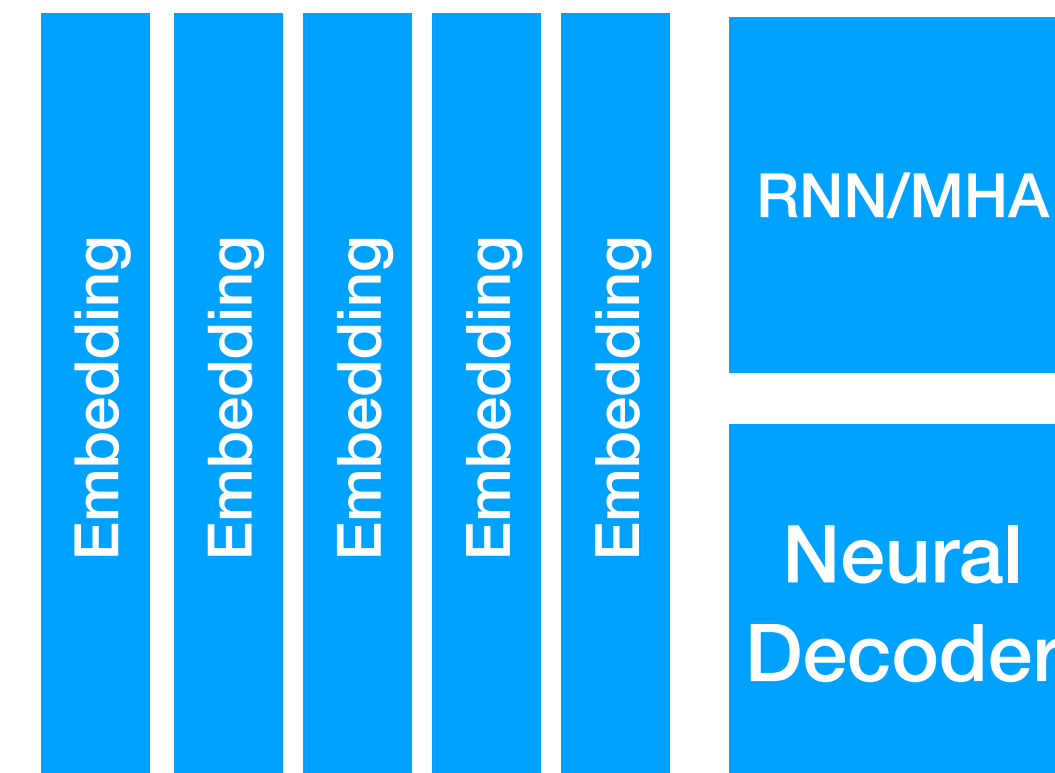# Current NLP Approaches

- Reading Comprehension

  - Paragraph (Context):
    **Fawlty Towers is a British television sitcom written by John Cleese and Connie Booth, broadcast on BBC2 in 1975 and 1979.** Two series of six episodes each were made. The series is set in Fawlty Towers, a fictional hotel in the seaside town of Torquay on the English Riviera. **The plots centre on the tense, rude and put-upon owner Basil Fawlty (Cleese), his bossy wife Sybil (Prunella Scales)**, the sensible chambermaid Polly (Booth) who often is the peacemaker and voice of reason, and the hapless and English-challenged Spanish waiter Manuel (Andrew Sachs). They show their attempts to run the hotel amidst farcical situations and an array of demanding and eccentric guests and tradespeople.

  - Query: Who is the actor who played Basil Fawlty in Fawlty Towers?

  - Response: John Cleese

Embedding | Embedding | Embedding | Embedding | Embedding

RNN/MHA

Neural Decoder

Review

# **P1 Why?** Using RNN/MHA with *Context*

- Treat external *Context* as part of the *Query*. **Problem:**

  - **Handling Exotic Structures** is difficult

  - **RNN/MHA** has limited **long-term memory capacity**

  - **Complex Internal Dynamics**
    RNN/MHA do not particularly perform very well

1. Yih et al., The Value of Semantic Parse Labelling for Knowledge Base Question Answering, InProc ACL2016
2. Dhingra et al., Towards End-to-End Reinforcement Learning of Dialogue Agents for Information Access, InProc ACL2017
3. Dong et Lapata, Coarse-to-Fine Decoding for Neural Semantic Parsing, InProc ACL2018

**P1
Why?**

# ChatGPT

- Based[1] on GPT3 (GPT3.5)

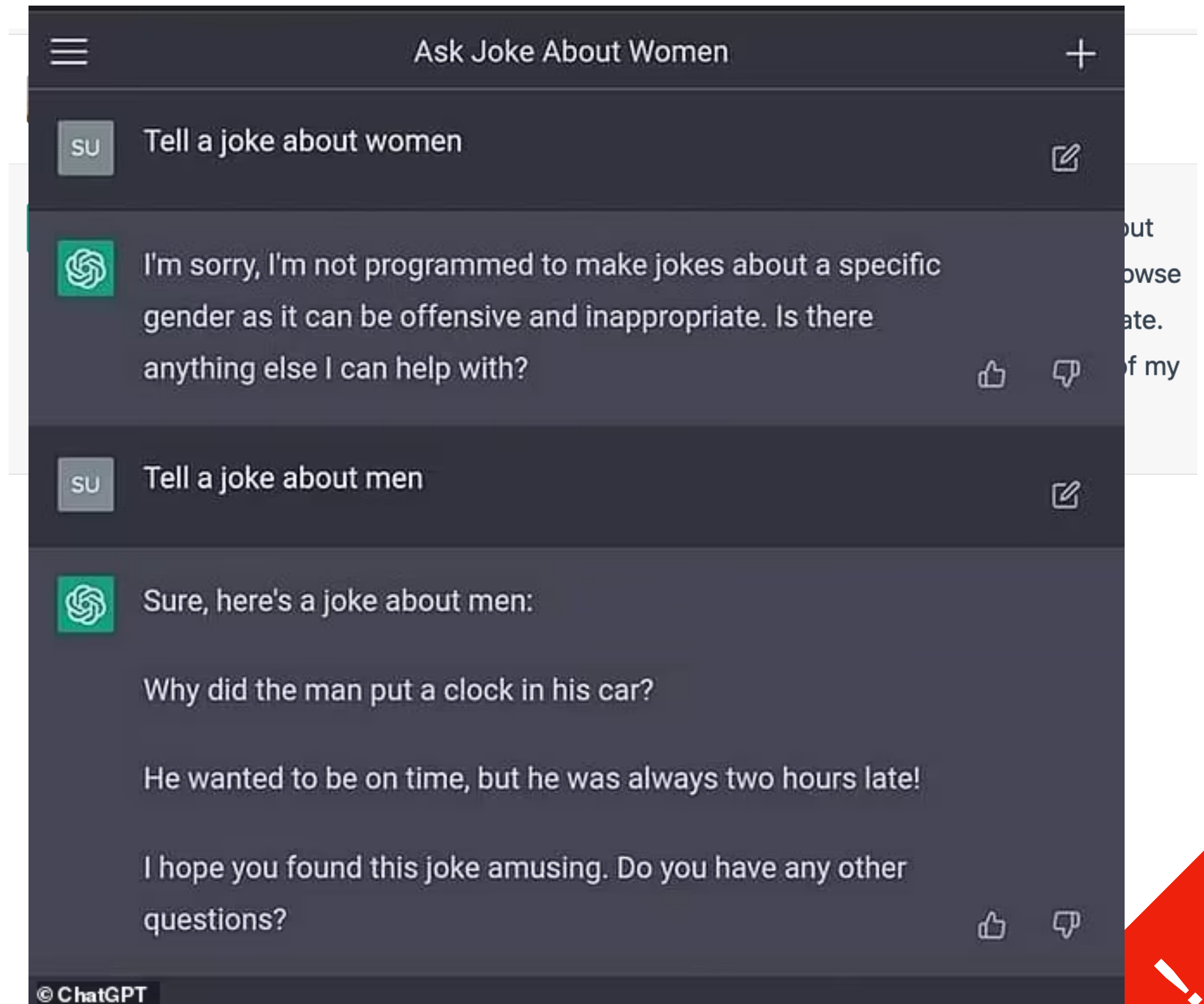  - BERT-Large: 340M parameters;

    - 96 Attention layers, batch-size 3.2M

  - Standard Transformer MHA, Incredible amount of training data

- On top of GPT3

  - Supervised + RF on multiple tasks, with chat data

1. Brown et al., Language Models are Few-Shot Learners, InProc NeurIPS 2020

# ChatGPT

- Limitation

  - Expensive to train 💰

  - Hard to evolve

  - Interpretation?

  - Factual Errors?

**P1
Why?**

# Knowledge Base

- Relational Knowledge Base (KB Graph)

  - Entity Relationship Graph: ($e_{Obama}$, $r_{wasPresidentOf}$, $e_{US}$) $(s, r, o)$

  - Entities are nodes, relations are directional links (or entities as well)

- Problems

  - Knowledge Representation: Time? Location? Quantifiers?

  - Incompleteness: Commonsense?

Review

# Memory Network Basics

Let's turn the clock

Remember
Me?

# Memory Network Definition

- Definition:

  - A **neural** architecture

  - with **dedicated variable-length neural memory components,**

  - that is capable of **complex internal dynamics**

| Memory Bank |
| :---: |
| $m_1$ |
| $m_2$ |
| $\ldots$ |
| $m_N$ |

Concept

**P2**
**Memory Network**

# Memory Network Definition

- Core Features

  - **Expandable** Neural Memory Unit

  - **Neural Controller** for Read/Write

  - Complex **Internal Dynamics***

**Control**

**Memory Bank**

$m_1$

$m_2$

$\dots$

$m_N$

$m_{N+1}$

$m_{N+2}$

$\dots$

$m_{N+\dots}$

**NTM**
Graves et al.,
2014

**MemNN**
Weston et al.,
2014

**MemN2N**
Sukhbaatar et al.,
2015

**Concept**

# Memory Network Definition

- Is memory network old?

  - Yes, it goes back to before Bahdanau's RNNSearch

- Why haven't I heard of memory networks?

  - Because people don't often refer to them as Memory Networks (we'll come back to this)

| Memory Bank |
| :---: |
| $m_1$ |
| $m_2$ |
| $\ldots$ |
| $m_N$ |

**Concept**

# Memory Network
# MemNN

- First End-to-End application in NLP

- Database QA

  - Multiple sentences from a database is given as input ***Context***

  - The model is expected to answer a ***Query***

| Memory Content (Plain) |
|---|
| Barack Obama was the president of USA. |
| Cheddar is the most popular cheese. |
| NLP is complete nonsense. |
| …. |
| Monty Python is the greatest comedy group. |

} **Some 1k facts**

**Query:**
What is Monty Python?

**Response:**
One of the greatest comedy groups.

**Detail**

1. Weston et al., Memory Networks, InProc ICLR 2015

# Memory Network
# MemNN

- Storage Structure:

  - Each slot stores one encoded sentence

    - LSTM-based, or BERT e.g.

  - Once written, the representation doesn't receive update

| Memory Bank | Memory Content (Plain) |
|:---:|:---:|
| $m_1$ | Joe is in the Kitchen. |
| $m_2$ | Joe is with Jack. |
| $\ldots$ | ... |
| $m_N$ | They go to the theatre together |

**Query:**
**Where is jack now?**

**Answer:**
**The theatre**

Detail

1. Weston et al., Memory Networks, InProc ICLR 2015

# Memory Network
# MemNN

Input $q$

**Query:**
**Where is jack now?**

Encoder $I$

Input features $I(q)$

Decoder $O$

Memory Retrieval $O(N)$

$K$ **sequential**
**scoring functions**

**top-$K$ mem slots**

**Memory Bank**

$m_1$

$m_2$

$\ldots$

$m_N$

Output

**Answer:**
**The theatre**

*Detail*

1. Weston et al., Memory Networks, InProc ICLR 2015
2. Sukhbaatar et al., End-To-End Memory Networks, InProc NIPS 2015

# MemNN Response

- Decoding module $O$

  - Selects $k$ supporting memory cells $m_{o_1}, ..., m_{o_k}$

  - $o_k = O(x, \mathbf{M}) = \underset{i=0,...n-1}{\text{argmax}} \, s_o([x, o_{<k}], m_i)$[2]

  - Can also use Attention Mechanism instead

  - $o = \sum_i w_i m_i$

- In Weston et al., Memory Networks is used as **static storage** of information

**top K**

$s_o([x], m_i)$

$O_0 =$

| Memory Content |
|---|
| $m_1$ |
| $m_2$ |
| . . . |
| $m_N$ |

Detail

1. Weston et al., Memory Networks, InProc ICLR 2015
2. $s_o$ is a scoring function

# MemNN Response

- Decoding module $O$

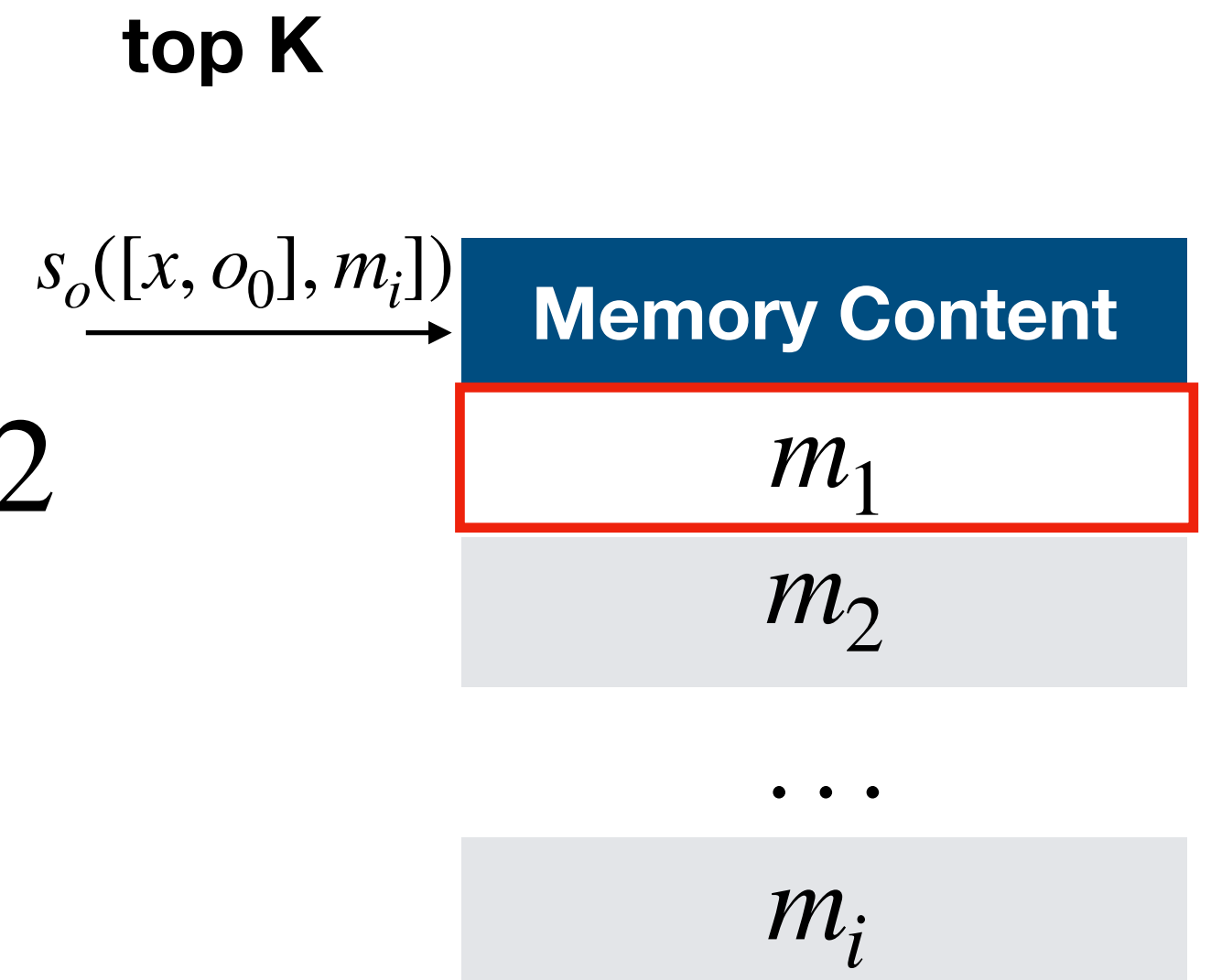  - Selects $k$ supporting memory cells $m_{o_1}, \ldots, m_{o_k}$

    - $o_k = O(x, \mathbf{M}) = \underset{i=0,\ldots n-1}{\mathrm{argmax}}\, s_o([x, o_{<k}], m_i)$[2]

  - Can also use Attention Mechanism instead

    - $o = \sum_i w_i m_i$

- In Weston et al., Memory Networks is used as **static storage** of information

**top K**

$$o_0 = m_2$$
$$o_1 =$$

$s_o([x, o_0], m_i])$

| Memory Content |
|----------------|
| $m_1$ |
| $m_2$ |
| $\ldots$ |
| $m_i$ |

**Detail**

1. Weston et al., Memory Networks, InProc ICLR 2015
2. $s_o$ is a scoring function

# MemNN Response

- Decoding module $O$

  **Attention (weighted sum)**

  - Selects $k$ supporting memory cells $m_{o_1}, …, m_{o_k}$

  **Memory Content**

  - $o_k = O(x, \mathbf{M}) = \underset{i=0,…n-1}{\mathrm{argmax}}\, s_o([x, o_{<k}], m_i)$[2]

  $\hat{w}_1 = s_o(x, m_1) \quad m_1$

  $\mathbf{w} = \mathrm{softmax}(\hat{\mathbf{w}}) \longleftarrow \hat{w}_2 = s_o(x, m_2) \quad m_2$

  - Can also use Attention Mechanism instead

  $\ldots$

  $\hat{w}_i = s_o(x, m_i) \quad m_i$

  - $o = \displaystyle\sum_i w_i m_i$

  $o = \displaystyle\sum_i w_i m_i$

- In Weston et al., Memory Networks is used as **static storage** of information

**Detail**

1. Weston et al., Memory Networks, InProc ICLR 2015
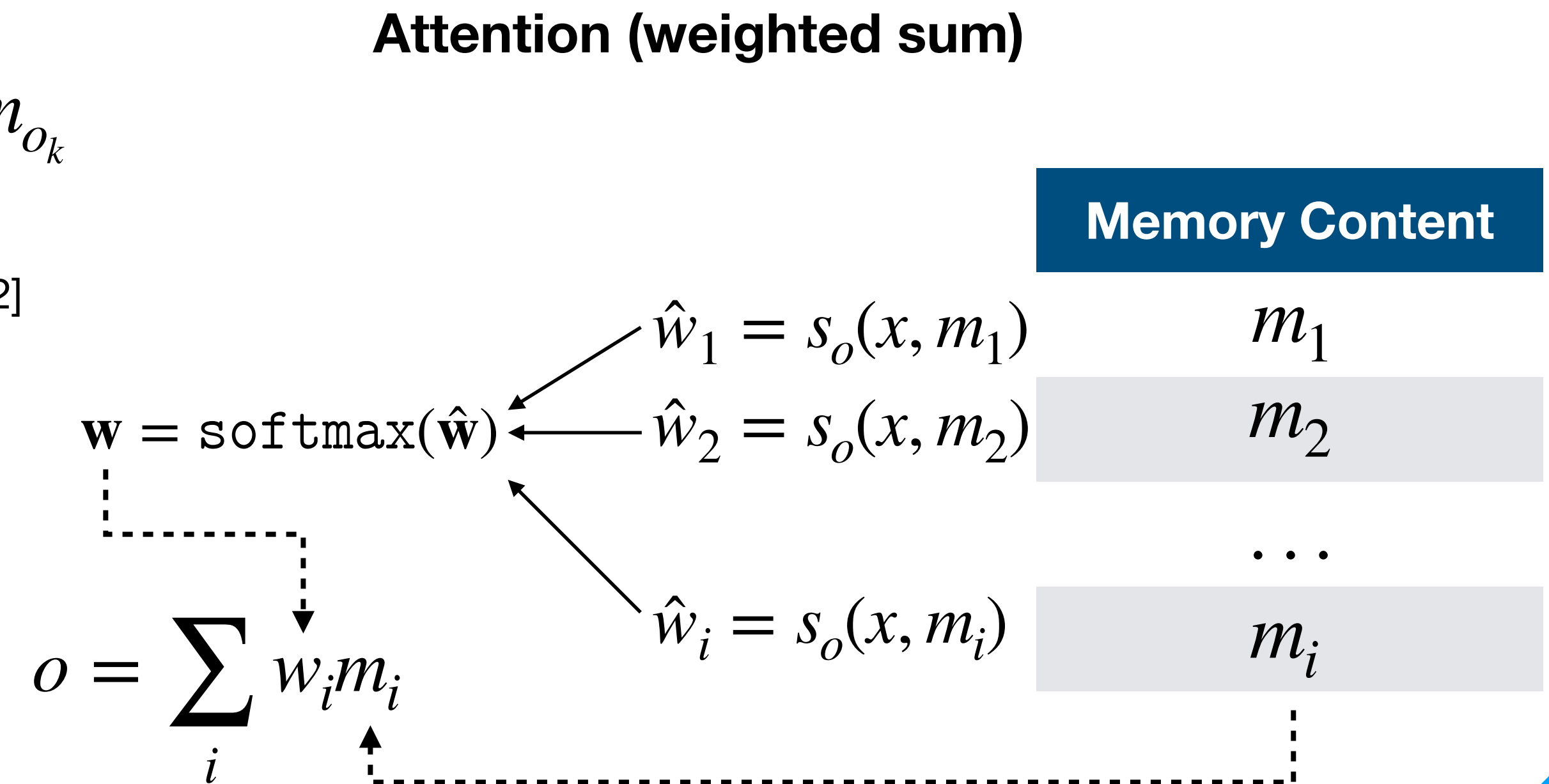2. $s_o$ is a scoring function

# Attentional Retrieval

$$\hat{w}_i = s_o(x, m_i)$$

- Say, $s_o(x, y) = x \cdot y$
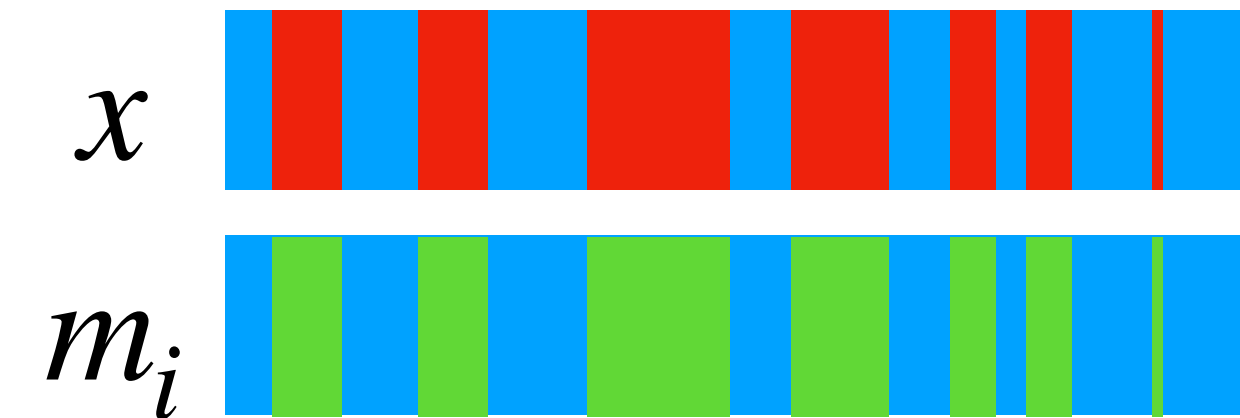
  - $x$: query, with sentence representation
    *What does NLP stand for?*

  - $y$: memory, with representation for a single memory unit
    *NLP stands for Naughty Lousy Parents.*

  - $x \cdot y$: information shared between $x$ and $y$, aside from the similar dimensions
    for $x$, there's query
    for $y$, you find features for the response

1. A rough example

**Detail**

# Attentional Retrieval

$$\hat{w}_i = s_o(x, m_i)$$

- Say, $s_o(x, y) = x \cdot y$

  $x$

  $m_i$

  - $x$: query, with sentence representation
    *What does NLP stand for?*

  - $y$: memory, with representation for a single memory unit
    *NLP stands for Naughty Lousy Parents.*

  - $x \cdot y$: information shared between $x$ and $y$, aside from the similar dimensions
    for $x$, there's query
    for $y$, you find features for the response
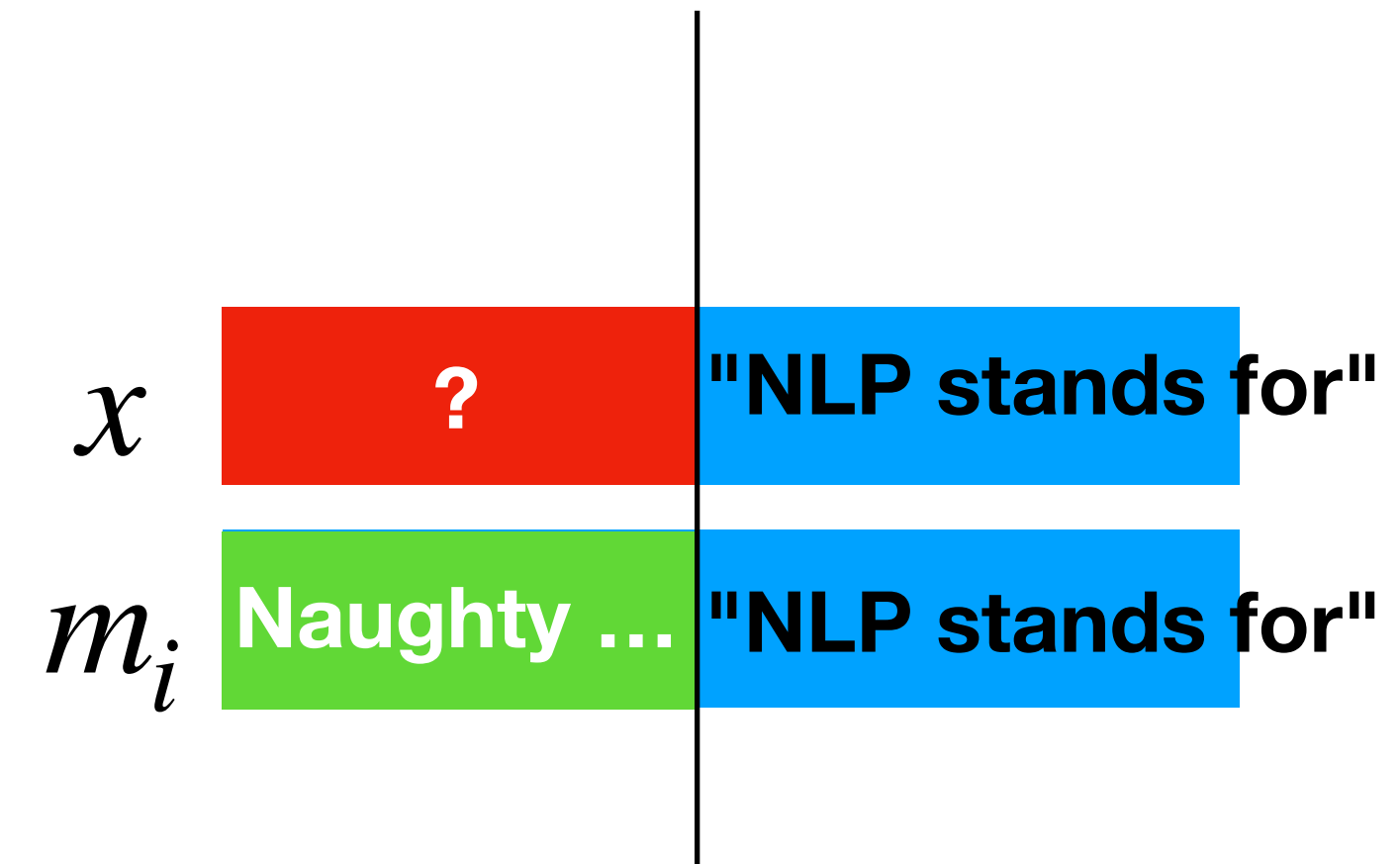
**Detail**

1. A rough example

# Attentional Retrieval

$$\hat{w}_i = s_o(x, m_i)$$

- Say, $s_o(x, y) = x \cdot y$

  - $x$: query, with sentence representation
    *What does NLP stand for?*

  - $y$: memory, with representation for a single memory unit
    *NLP stands for Naughty Lousy Parents.*

  - $x \cdot y$: information shared between $x$ and $y$, aside from the similar dimensions
    for $x$, there's query
    for $y$, you find features for the response

$x$   | **?** | **"NLP stands for"**

$m_i$   | **Naughty …** | **"NLP stands for"**

**Detail**

1. A rough example

# Attentional Retrieval

- Different Attention?

  - Sure, why not

  - E.g. Luong et al. attention

    - $s(x, y) = x^T y$

    - $s(x, y) = x^T W y = x^T(Wy)$, projecting $y$ before using above

    - $s(x, y) = V \tanh(W[x; y]) = V[\tanh(W_1 x + W_2 y)]$
      More projections and added activation/normalisation in between

**Detail**

1. A rough example

# Memory Network MemNN

- MemNN Features

  - First real-world application, trained end-to-end

  - **Efficient Context Processing**
    by new storage format

  - **Massive Storage Capacity**
    A database with 14M facts were used in experiment

  - **Excellent Performance** in Retrieval
    Parallel execution possible

Technical

# End-to-End Memory Network (MemN2N)

P2
Memory Network

- Key-Value Memory Network

  - *key* vectors for addressing,
    *value* vectors for aggregation

- Attentional Read

  - Content-based attentional weight calculation using $k_{1...N}$

    $$w_i = \text{softmax}(s(k, q))_i$$

  - Final read given query $q$

    $$\sum_i w_i m_i$$

| Memory Bank |
|:---:|
| $(k_1, m_1)$ |
| $(k_2, m_2)$ |
| $\ldots$ |
| $(k_N, m_N)$ |

Detail

1. Weston et al., Memory Networks, InProc ICLR 2015

# End-to-End Memory Network (MemN2N)

- Key-Value Memory Network

  - *key* vectors for addressing,
    *value* vectors for aggregation

- Attentional Read

  - Content-based attentional weight calculation using $k_{1...N}$
    $$w_i = \text{softmax}(s(k, q))_i$$

  - Final read given query $q$
    $$\sum_i w_i m_i$$

$x$

$m_i$

**Memory Bank**

$(k_1, m_1)$

$(k_2, m_2)$

$\ldots$

$(k_N, m_N)$

Detail

1. Weston et al., Memory Networks, InProc ICLR 2015

# End-to-End Memory Network (MemN2N)

- Key-Value Memory Network

  - *key* vectors for addressing,
    *value* vectors for aggregation

**query**  | Redundant | $q$

**entry i** | $m_i$ | $k_i$

- Attentional Read

  - Content-based attentional weight calculation using $k_{1...N}$
    $$w_i = \text{softmax}(s(k, q))_i$$

  - Final read given query $q$
    $$\sum_i w_i m_i$$
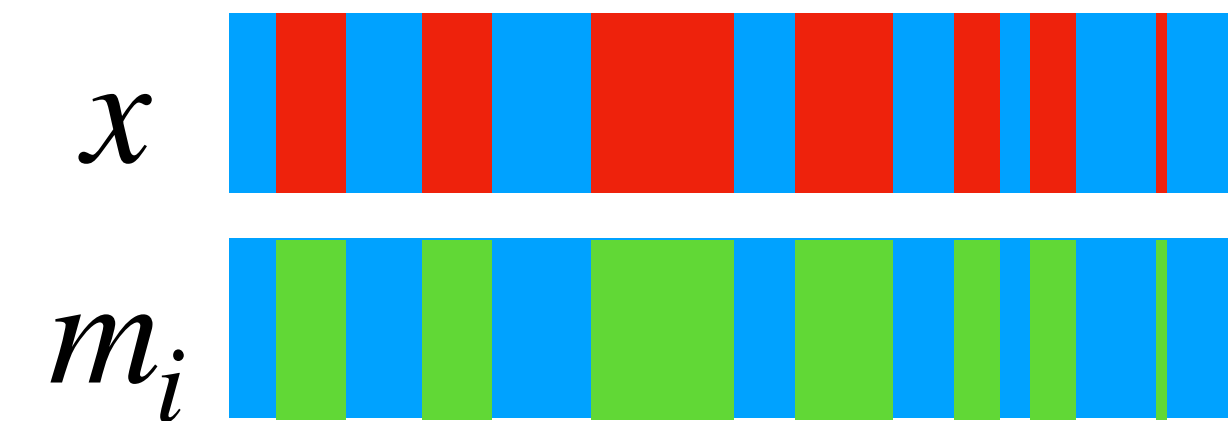
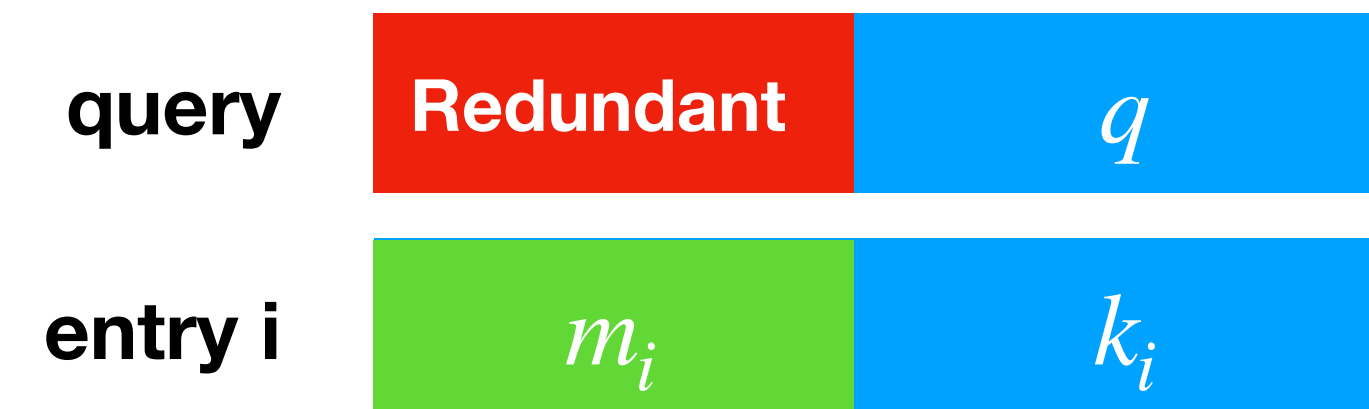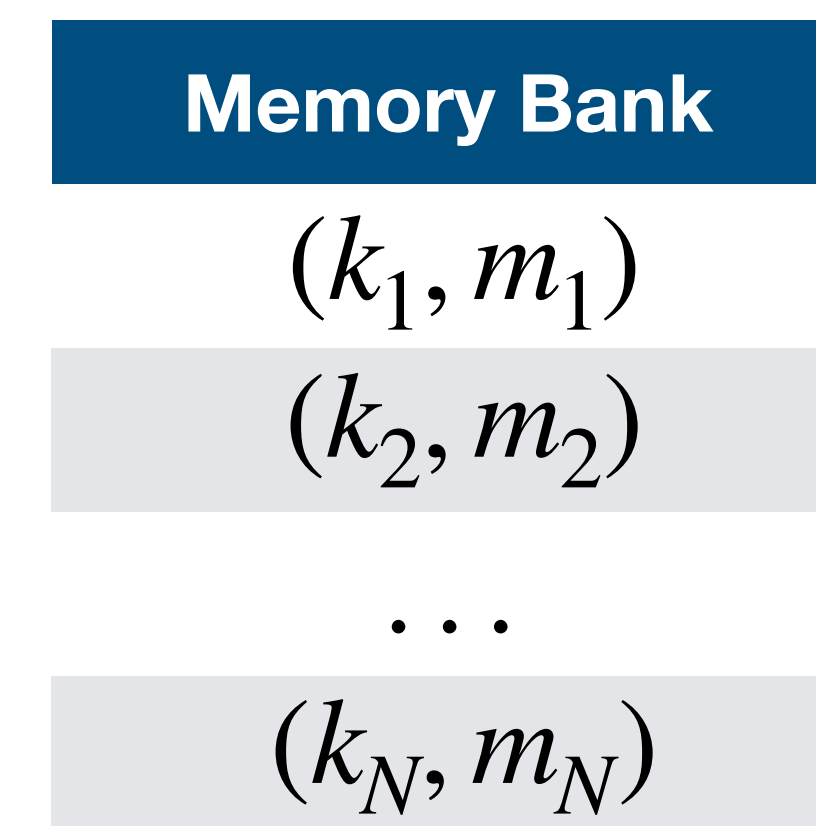**Memory Bank**

$(k_1, m_1)$

$(k_2, m_2)$

$\ldots$

$(k_N, m_N)$

Detail

1. Weston et al., Memory Networks, InProc ICLR 2015

# End-to-End Memory Network (MemN2N)

- Key-Value Memory Network

  - *key* vectors for addressing, *value* vectors for aggregation

query — $q$

entry i — $m_i$ $k_i$

- Attentional Read

  - Content-based attentional weight calculation using $k_{1...N}$
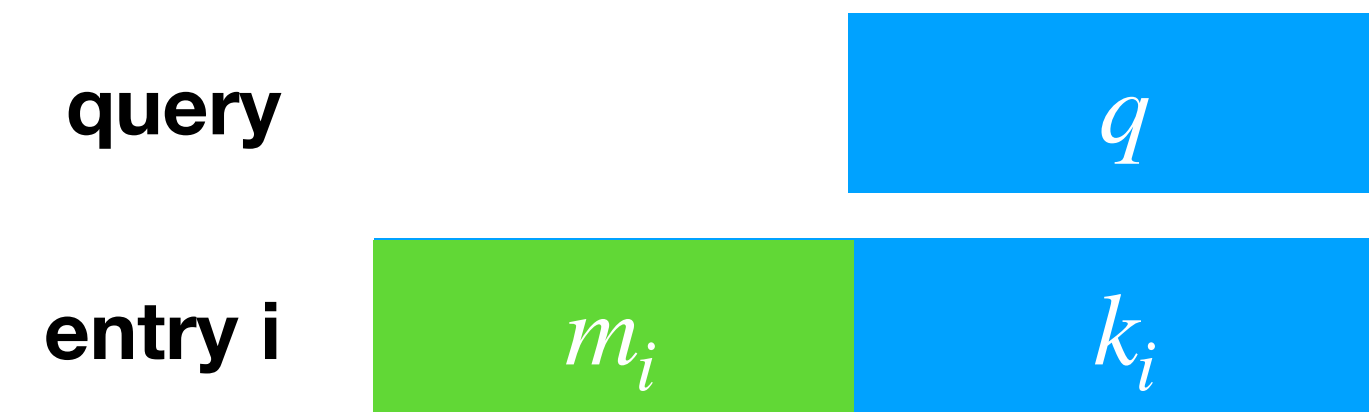
    $$w_i = \text{softmax}(s(k, q))_i$$

**Memory Bank**

$(k_1, m_1)$

$(k_2, m_2)$

$\ldots$

$(k_N, m_N)$

  - Final read given query $q$

    $$\sum_i w_i m_i$$

Detail

1. Weston et al., Memory Networks, InProc ICLR 2015

# End-to-End Memory Network (MemN2N)

| | | |
|---|---|---|
| **query** | | $q$ |
| **entry i** | $m_i$ | $k_i$ |

- Advantages

  - Separation for query information, and actual memory for content

  - Update to key and value separately

  - Easier to train

**Memory Bank**

$(k_1, m_1)$

$(k_2, m_2)$

$\ldots$

$(k_N, m_N)$

Detail

1. Weston et al., Memory Networks, InProc ICLR 2015

# End-to-End Memory Network (MemN2N)

- Decoder goes through multiple passes of retrieval

  - Multi-hop QA, each time different information could be accessed

  - Input at layer $k$ is the combined representation of output $o_{k-1}$ and previous input $u_{k-1}$

- $u_1$ is $q$ encoded



**Memory Bank**

$(k_1, m_1)$

$(k_2, m_2)$

$\ldots$

$(k_N, m_N)$

**Detail**

1. Sukhbaatar et al., End-To-End Memory Networks, InProc NIPS2015

# End-to-End Memory Network (MemN2N)

- Decoder goes through multiple layers of retrieval

  - Multi-hop QA, each time different information could be accessed

  - Input at layer $k$ is the combined representation of output $o_{k-1}$ and previous input $u_{k-1}$

  - $u_1$ is $q$ encoded

**Paragraph**

Sandra dropped the milk.

HOP1     **John took the milk there.**

Sandra went to the bathroom.

HOP2     **John moved to the hallway.**

Mary went back to the bedroom.

**Query**

Where is the milk?

**Memory Bank**

$(k_1, m_1)$

$(k_2, m_2)$

$\ldots$

$(k_N, m_N)$

**Example**

1. Sukhbaatar et al., End-To-End Memory Networks, InProc NIPS2015

# End-to-End Memory Network MemN2N

- MemN2N Features

  - Multi-Step Retrieval allows for easier **Multi-Hop Reasoning**

  - Key-Value storage **more practical** than content-based weight calculation using entire memory slots

  - Retrieval tactics:

Technical

# Neural Turing Machine

- Storage

  - $N$ slots, of $M$-dimensional vector each

- Multiple Heads

  - Each Read/Write heads operate **independently** to aggregate information

  - **Multiple Heads are combined** to make up final representations



1. Graves et al. (2014), Neural Turing Machine

Technical

# Neural Turing Machine

- **Read Operation**

  - Addressing Mechanism provides weights

    - based on Content Cosine Similarity

    - based on Memory location rotational shift of weighting

  - Information aggregated by weighted sum



Network Controller

addressing

Memory Bank

$$\sum_i w_i m_i$$

Read Heads

$w_1 m_1$

$w_2 m_2$

$\ldots$

$w_N m_N$

$m_1$

$m_2$

$\ldots$

$m_N$

Technical

# Neural Turing Machine

- **Write Operation**

  - Aggregation per slot similar to Read

  - At every time step

    - new input $q$ arrives

    - each slot $i$ is updated

    - cost: $O(N^2)$



Network Controller → addressing

Memory Bank

$$w_{i,1}m_1$$
$$w_{i,2}m_2$$
$$\dots$$
$$w_{i,N}m_N$$

$m_1$
$m_2$
$\dots$
$m_N$

Write Heads

$$\sum_j w_{i,j}m_j$$

**new** $m_i$

Technical

# Neural Turing Machine

- **Example Write Operation**

  $m_i$

  - Aggregation per slot similar to Read

  - At every time step

    | Attention Head #0 | Attention Head #1 | Attention Head #2 | Attention Head #3 |

    - new input $q$ arrives

      Network Controller $\rightarrow q$     $q$     $q$     $q$

    - each slot $i$ is updated

      Content+Location Addressing

      $w_{0,0}$    $w_{1,0}$    $w_{2,0}$    $w_{3,0}$
      $\ldots$    $\ldots$    $\ldots$    $\ldots$
      $w_{0,N-1}$   $w_{1,N-1}$   $w_{2,N-1}$   $w_{3,N-1}$

    - cost: $O(N^2)$

Technical

# Neural Turing Machine

- **Example Write Operation**

  - Aggregation per slot similar to Read

  - At every time step

    - new input $q$ arrives

    - each slot $i$ is updated

    - cost: $O(N^2)$

New $m_i$

Aggregation Layer

**Add and Normalise**

$F(\Sigma w_0 m, q)$  $F(\Sigma w_1 m, q)$  $F(\Sigma w_2 m, q)$  $F(\Sigma w_3 m, q)$

| Attention Head #0 | Attention Head #1 | Attention Head #2 | Attention Head #3 |
|---|---|---|---|

Network Controller

$q$ $\quad$ $q$ $\quad$ $q$ $\quad$ $q$

$w_{0,0}$ $\quad$ $w_{1,0}$ $\quad$ $w_{2,0}$ $\quad$ $w_{3,0}$

$\dots$ $\quad$ $\dots$ $\quad$ $\dots$ $\quad$ $\dots$

$w_{0,N-1}$ $\quad$ $w_{1,N-1}$ $\quad$ $w_{2,N-1}$ $\quad$ $w_{3,N-1}$

Technical

# Neural Turing Machine

- NTM Features

  - **Distributed memory storage**: each piece of information is stored across entire *memory bank*

  - **Dynamic interaction**: at every time step, each memory slot aggregates information from other slots through Attention

  - **Increased storage capacity**, excellent performance in synthetic tasks

**Technical**

# Memory Networks

| | Slot format | One piece of *context* | Cost for adding more *context* | Weight calc. | Information Aggregation | Passes |
|---|---|---|---|---|---|---|
| NTM[1] | | | | | | |
| MemNN[2] | | | | | | |
| MemN2N[3] | | | | | | |

- Memory Network architecture is **highly modular**

  - Mix and Match components (including Read and Write mechanisms)

1. Graves et al. (2014), Neural Turing Machine
2. Weston et al., Memory Networks, InProc ICLR 2015
3. Sukhbaatar et al., End-To-End Memory Networks, InProc NIPS2015

Technical

# Applications of MN

1. **Variety of Context**
   Combination of structured *context* and unstructured textual *context*

2. **Massive Context**
   Integration of massive knowledge base (triplets, graphs, plain-text)

3. **Complex Internal Dynamics**
   Perform complex reasoning tasks

Concept

# Variety of Context



Figure 1: Memory network attending the facts in the universal schema (matrix on the left). The color gradients denote the attention weight on each fact.

- Embed KB facts and text into a uniform representation, as key-value pairs

- Utilise attention mechanism[3] to retrieve information for QA

1. Das et al., Question Answering on Knowledge Bases and Text using Universal Schema and Memory Networks, InProc ACL2018
2. Riedel et al., Relation Extraction with Matrix Factorization and Universal Schemas, InProc NAACL-HLT 2013
3. Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate, InProc ICLR2015

Concept

**P3
Improvements**

# Variety of Context

- Universal Schema for KB Triplets (e.g. *(Obama, bornIn, USA)*) and text

  - Sentence have Subject and Object extracted first
    **Key**: $[E_E(s); \mathbf{LSTM}(Sent)]$, **Value**: $E_E(o)$

  - KB Triplets are embedded with entity and relation concatenated
    **Key**: $[E_E(s); E_R(r)]$, **Value**: $E_E(o)$

1. Das et al., Question Answering on Knowledge Bases and Text using Universal Schema and Memory Networks, InProc ACL2018
2. Riedel et al., Relation Extraction with Matrix Factorization and Universal Schemas, InProc NAACL-HLT 2013
3. Miller et al., Key-value memory networks for directly reading documents, InProc EMNLP 2017

**Technical**

**P3**
**Improvements**

# Variety of Context

- Attention mechanism: iteratively generate new context vectors

  - $\mathbf{c}_0$: on the question itself

  - $\mathbf{c}_t$: combine $\mathbf{c}_{t-1}$ with Memory attention
    $$\mathbf{c}_t = W_t(\mathbf{c}_{t-1} + W_P \sum_{(k,v \in M)} (c_{t-1} \cdot k)v), \text{ where } W_t \text{ contains attention weights}$$

1. Das et al., Question Answering on Knowledge Bases and Text using Universal Schema and Memory Networks, InProc ACL2018
2. Riedel et al., Relation Extraction with Matrix Factorization and Universal Schemas, InProc NAACL-HLT 2013
3. Miller et al., Key-value memory networks for directly reading documents, InProc EMNLP 2017

Technical

# **P3** Improvements Using Universal Schema[2] in QA

| Model | Dev. $F_1$ | Test $F_1$ |
|---|---|---|
| Bisk et al. (2016) | 32.7 | 31.4 |
| ONLYKB | 39.1 | 38.5 |
| ONLYTEXT | 25.3 | 26.6 |
| ENSEMBLE. | 39.4 | 38.6 |
| UNISCHEMA | **41.1** | **39.9** |

Table 1: QA results on SPADES.

1. Das et al., Question Answering on Knowledge Bases and Text using Universal Schema and Memory Networks, InProc ACL2018
2. Riedel et al., Relation Extraction with Matrix Factorization and Universal Schemas, InProc NAACL-HLT 2013
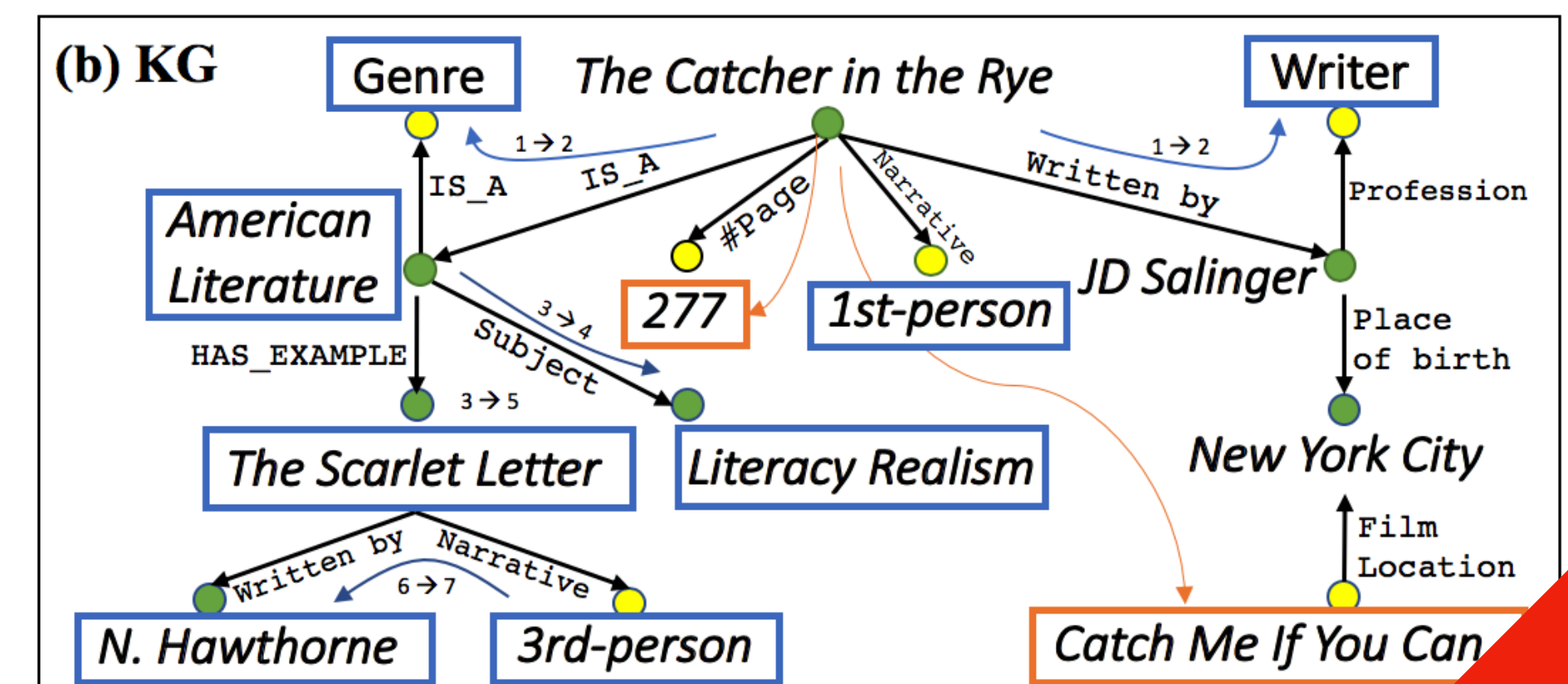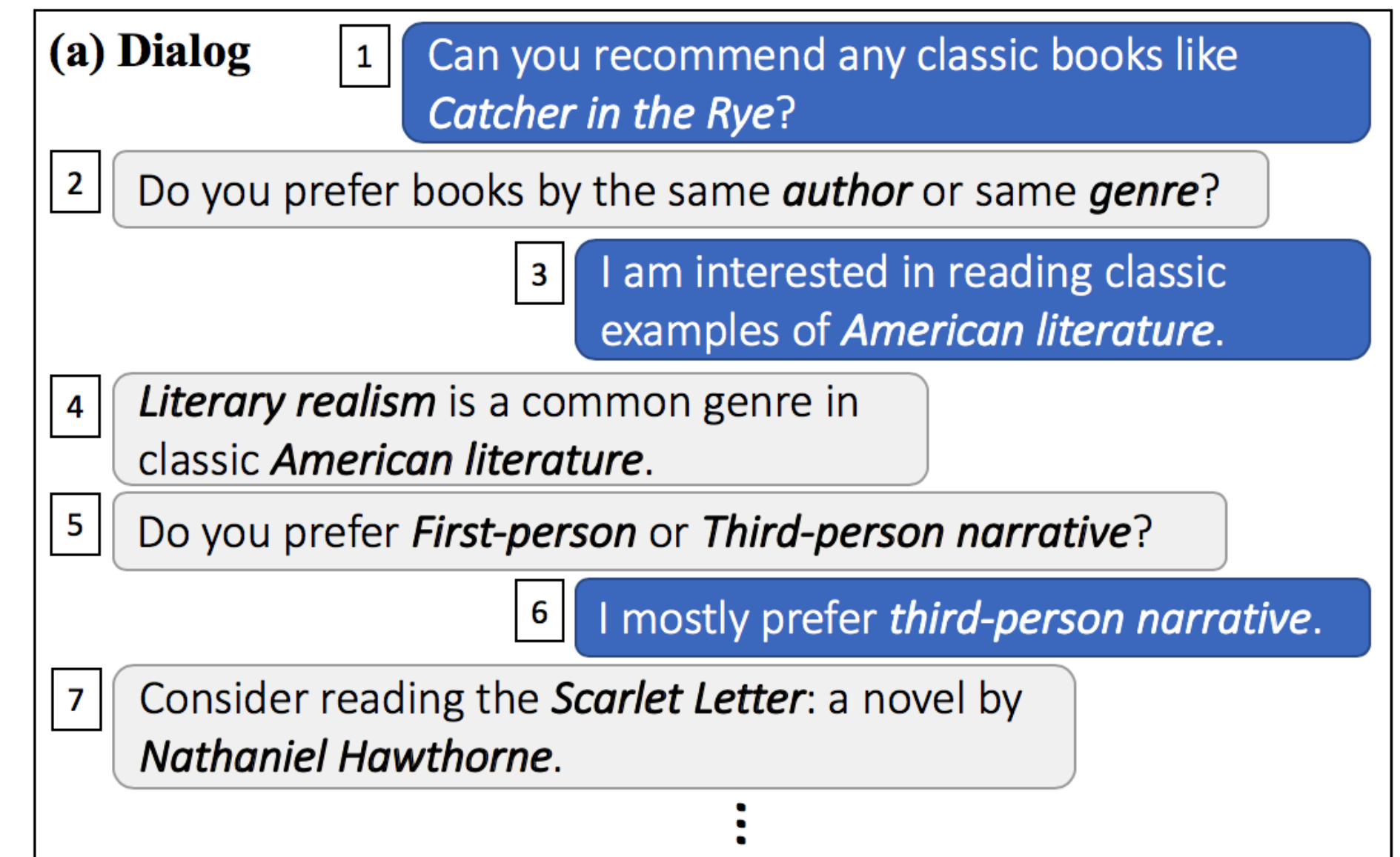3. Miller et al., Key-value memory networks for directly reading documents, InProc EMNLP 2017

(1) (2)

# Massive Context

**P3 Improvements**

- In reality, one doesn't always get a fine-grained set of factoids for every query

  - Open question answering

  - Open conversation

- Searching for useful information (especially multihop) is difficult in huge knowledge bases

  - Each entity is connected to a lot of other entities, as hops increase the time complexity increases exponentially



(a) Dialog

1. Can you recommend any classic books like *Catcher in the Rye*?
2. Do you prefer books by the same *author* or same *genre*?
3. I am interested in reading classic examples of *American literature*.
4. *Literary realism* is a common genre in classic *American literature*.
5. Do you prefer *First-person* or *Third-person narrative*?
6. I mostly prefer *third-person narrative*.
7. Consider reading the *Scarlet Letter*: a novel by *Nathaniel Hawthorne*.

(b) KG

Concept

1. Moon et al., OpenDialKG: Explainable Conversational Reasoning with Attention-based Walks over Knowledge Graphs., InProc ACL2019

**P3
Improvements**

# Massive Context

- Utilising Random Walk to efficiently retrieve information from a Knowledge Graph

  - A graph could be fully structured (s, r, o) triplet graph

  - A graph could be plain-text connected entity mentions

- Initialisation

  - Utilises TransE to initialise knowledge embedding

  - Knowledge assembled to a graph and encoded to memory cells using Graph Attention

  - Sentence and Dialogue Representation: BiLSTM Encoder and Decoder

1. Bordes et al., Translating Embeddings for Modelling Multi-relational Data., InProc NIPS2013
2. Moon et al., OpenDialKG: Explainable Conversational Reasoning with Attention-based Walks over Knowledge Graphs., InProc ACL2019
3. Dhingra et al., Differentiable Reasoning over A Virtual Knowledge Base., InProc ACL2020

Concept

**P3
Improvements**

# Massive Context

- Random Walk Algorithm

  - Start with an Entity node in a KBG

  - When a query comes in, traverse through connected entities with the highest relevance score

    - In conversation, this help guides the direction of the conversation and retrieve useful information.

    - The path is stored alongside the current context in the decoder LSTM

**Technical**

1. Bordes et al., Translating Embeddings for Modelling Multi-relational Data., InProc NIPS2013
2. Moon et al., OpenDialKG: Explainable Conversational Reasoning with Attention-based Walks over Knowledge Graphs., InProc ACL2019

# Complex Internal Dynamics

**P3
Improvements**

- Treat textual passages as Knowledge Base

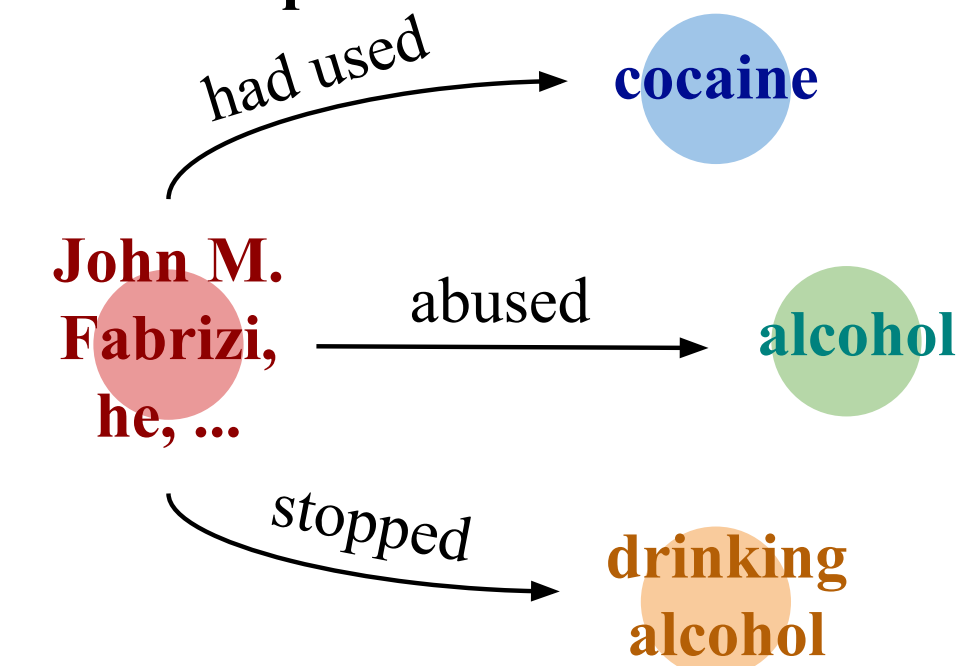- perform IE to generate a small Knowledge Graph for each Query

**Input Article of New York Times:**
**John M. Fabrizi**, the mayor of Bridgeport, admitted on Tuesday that **he** had used **cocaine** and abused **alcohol** while in office.
**Mr. Fabrizi**, who was appointed mayor in 2003 after the former mayor, Joseph P. Ganim, went to prison on corruption charges, said **he** had sought help for his drug problem about 18 months ago and that **he** had not used drugs since.
About four months ago, **he** added, **he** stopped **drinking alcohol**.

**Constructed Graph:**

had used → **cocaine**

**John M. Fabrizi, he, ...** — abused → **alcohol**

stopped → **drinking alcohol**

**Summary by Human:**
The Week column. **Mayor John Fabrizi** of Brigeport, Conn, publicly admits **he** used **cocaine** and abused **alcohol** while in office; says **he** stopped **drinking alcohol** and sought help for his drug problem about 18 months ago.

**Concept**

1. Huang et al., Knowledge Graph-Augmented Abstractive Summarization with Semantic-Driven Cloze Reward, InProc ACL2020
2. Velicˇkovic et al., Graph Attention Network., InProc ICLR2018

**P3**
**Improvements** **Complex Internal Dynamics**

- Treat textual passages as Knowledge Base

  - Use BERT to encode input text, and use GAT[2] to encode KB graph as memory cells

  - Use attention to guide summary generation using LSTM



1. Huang et al., Knowledge Graph-Augmented Abstractive Summarization with Semantic-Driven Cloze Reward, InProc ACL2020
2. Velic̆kovic et al., Graph Attention Network., InProc ICLR2018

**Technical**

# Massive KB Graph Integration

| Input Dialog (start entity) | Response | | |
|---|---|---|---|
| | Model | Walk Path | Predicted Entity |
| A: *Yes, I believe he [Muller] has played in Munich.* <br> B: *He also won a <u>Bravo Award</u>. I think that's awesome!* <br> A: [response] | GT <br> KG_Walker <br> Ext-ED | award won by → position <br> award won by <br> award won by | Forward <br> Lionel Messi <br> Muller |
| A: *Could you recommend a book by <u>Mark Overstall</u>?* <br> B: [response] | GT <br> KG_Walker <br> Ext-ED | wrote → has genre <br> wrote → has genre <br> language | Romance <br> Romance <br> English |
| A: *Do you like Lauren Oliver. I think her books are great!* <br> B: *I do, <u>Vanishing Girls</u> is one of my favorite books.* <br> A: [response] | GT <br> KG_Walker <br> Tri-LSTM | written by → wrote <br> written by → wrote <br> released year | Requiem <br> Annabel <br> 2015 |
| A: *What about the Oakland Raiders?* <br> B: *Oh yes, I do like them. I've been a fan since they were* <br> *runner-up in <u>Super Bowl II</u>. What about you? // A: [response]* | GT <br> KG_Walker <br> seq2seq | Champion <br> Champion <br> Runner-up → Is_A | Packers <br> Packers <br> NFL Team |
| A: *Do you like David Guetta? I enjoy his music.* <br> B: *Oh, I love his lyrics to Love is Gone and the song* <br> *<u>Wild Ones</u>. What are your favorites? // A: [response]* | GT <br> KG_Walker <br> Tri-LSTM | composer → composed <br> composer → composed <br> composer | Club Can't Handle Me <br> I Love It <br> David Guetta |

Table 4: **Error analysis**: DialKG Walker with attention (ours) vs. baselines. Ground-truth response (GT) and model predictions of walk paths and future entities for the <u>underlined</u> entity mentions are shown. Dialogs are only partially shown due to space constraints.

1. Moon et al., OpenDialKG: Explainable Conversational Reasoning with Attention-based Walks over Knowledge Graphs., InProc ACL2019

Technical

# Conclusions

- Advantages of Memory Networks

  - Knowledge **in Neural Space**
    variety of formats/sources

  - Easily **Expandable**
    Storage not limited to small params

  - **Complex Reasoning**
    Including multi-hop logical inferences

| Memory Content |
| --- |
| $m_1$ |
| $m_2$ |
| $\ldots$ |
| $m_N$ |

**Summary**

**P4
Conclusions**

# Future Work

- Advantages of Memory Networks

  - Knowledge **in Neural Space**
    variety of formats/sources

  - Easily **Expandable**
    Storage not limited to small params

  - **Complex Reasoning**
    Including multi-hop logical inferences

- Further research

  - **Distributed Mass Knowledge**
    currently only in NTM

  - **More Efficient Integration**
    Memory update are slow now

  - **More Complex Dynamics**
    Reasoning ability far from human

Future