

Commonsense Knowledge Aware Conversation Generation with Graph Attention

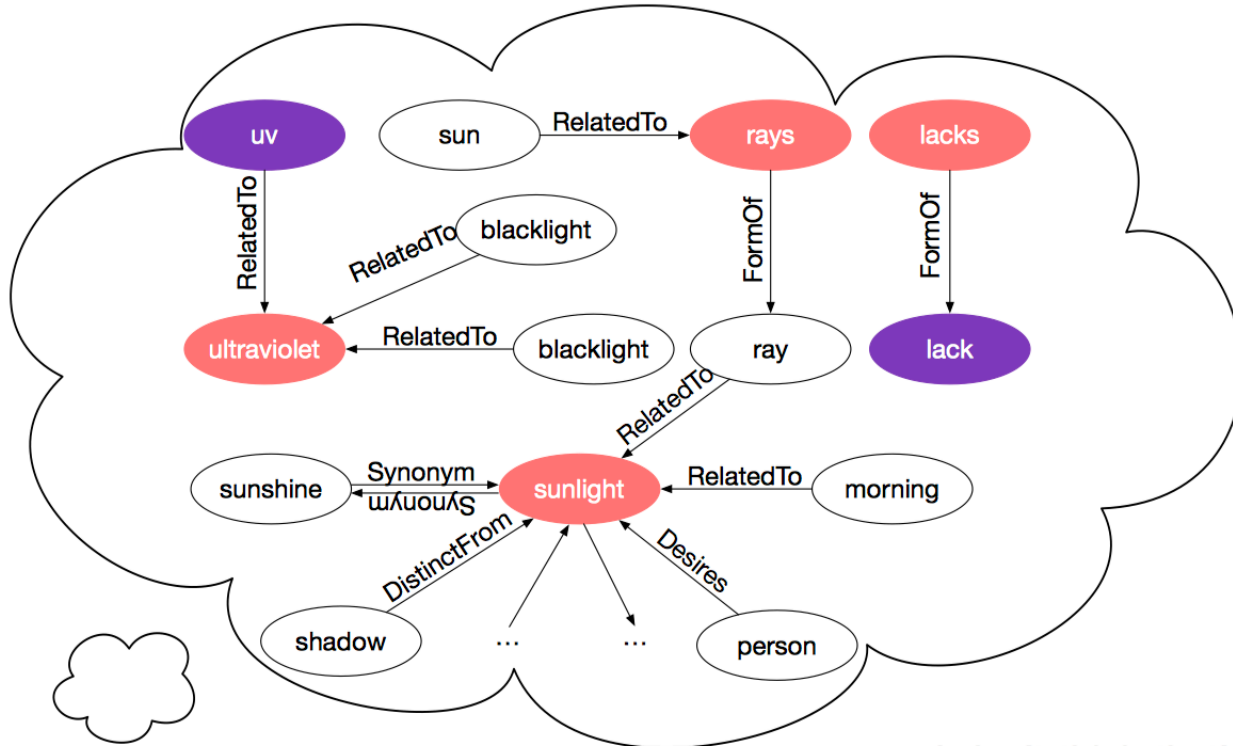
IJCAI 2018

Hao Zhou, Tom Young, **Minlie Huang**, Haizhou Zhao,
Jingfang Xu, Xiaoyan Zhu

Outline

1. Overview
2. Problem
3. Motivation
4. Encoder-Decoder
5. Task Definition
6. Model-CCM
7. Experiment
8. Conclusion

Overview



Moonlight **lacks** the **ultraviolet rays** of **sunlight**. → I don't think that's a **lack** of **uv**.

Moonlight lacks the ultraviolet rays of sunlight. → I'm not sure what you're saying.

with/without considering commonsense knowledge

Problem

- **Conversational data**
 - Tend to generate **generic responses**, which are unable to respond **appropriately** and **informatively**
- **Unstructured texts or domain-specific knowledge triples**
 - highly dependent on the **quality of unstructured texts**
 - **limited** by the small-scale, domain-specific knowledge
 - make use of knowledge triples (entities) **separately** and **independently**, instead of treating knowledge triples as a whole in a graph

Motivation

- **Large-scale commonsense knowledge**
 - understand the **background information** of a given post
 - facilitate response generation with such knowledge
- **Attention mechanism**
 - **Static graph** attention mechanism (*for encode*)
 - Encodes the **retrieved graphs** for a post to **augment** the semantic representation of the post
 - **Dynamic graph** attention mechanism (*for decode*)
 - Reads the knowledge graphs and the triples in each graph for better **response generation**

Encoder-Decoder

- **Encoder**

- Post: $X = x_1 x_2 \cdots x_n$

- Hidden representations: $h_t = \mathbf{GRU}(h_{t-1}, e(x_t))$,

- **Decoder**

- Hidden representations: $s_t = \mathbf{GRU}(s_{t-1}, [c_{t-1}; e(y_{t-1})])$,

- **Probability distribution:**

$$\begin{aligned} y_t \sim \mathbf{o}_t &= P(y_t \mid y_{<t}, \mathbf{c}_t) \\ &= \text{softmax}(\mathbf{W}_o \mathbf{s}_t). \end{aligned}$$

Task Definition

- **Given**

- Post: $X = x_1x_2 \cdots x_n$

- Commonsense knowledge graphs: $G = \{g_1, g_2, \cdots, g_{N_G}\}$

- *Each word corresponds to a graph in G*

- *Each graph consists of a set of triples*

$$g_i = \{\tau_1, \tau_2, \cdots, \tau_{N_{g_i}}\}$$

- *Each triple (head entity, relation, tail entity) is denoted as $\tau = (h, r, t)$*

- *knowledge triple τ is represented by*

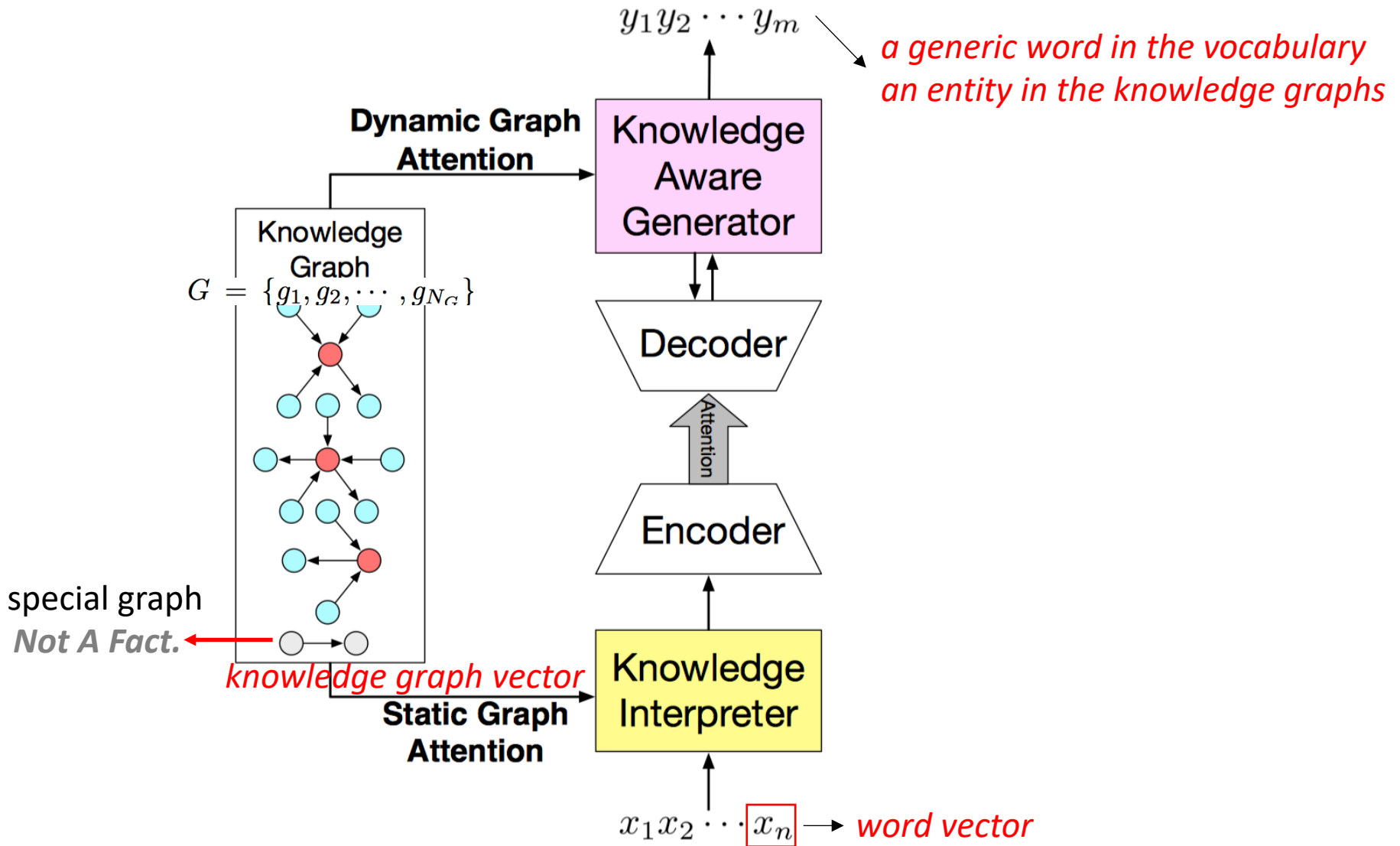
$$\mathbf{k} = (\mathbf{h}, \mathbf{r}, \mathbf{t}) = \mathbf{MLP}(\text{TransE}(\hat{h}, r, t))$$

- **Goal**

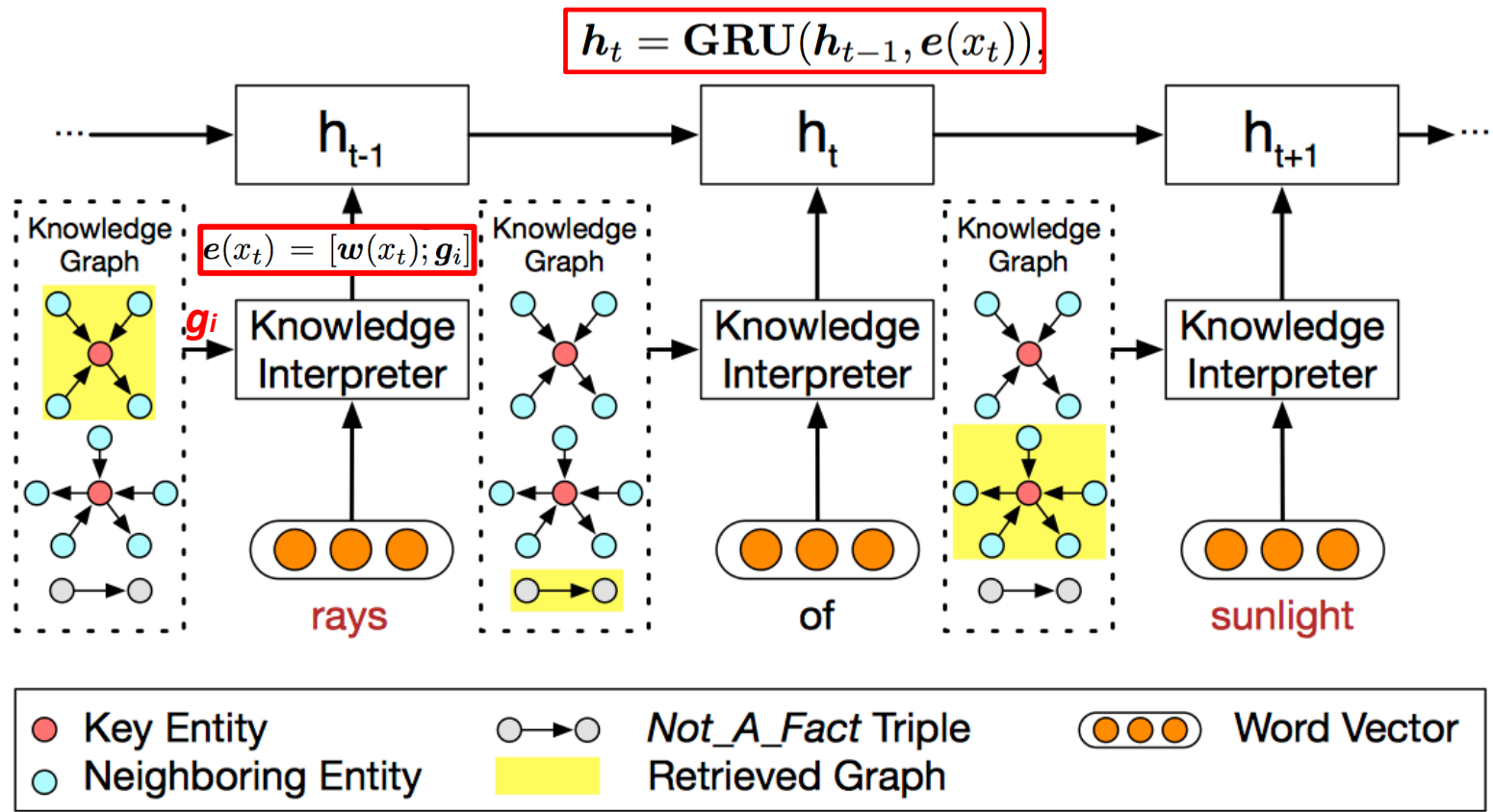
- response: $Y = y_1y_2 \cdots y_m$

- **Probability** : $P(Y|X, G) = \prod_{t=1}^m P(y_t|y_{<t}, X, G)$.

Model-CCM



Knowledge Interpreter



Static Graph Attention

- **Generate a representation for a retrieved knowledge graph**
 - not only all **nodes** in a graph
 - but also **relations** between nodes
- **Input:** $K(g_i) = \{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_{N_{g_i}}\}$
 $\mathbf{k} = (\mathbf{h}, \mathbf{r}, \mathbf{t}) = \text{MLP}(\text{TransE}(\hat{h}, r, t))$
- **Output:** **graph vector** \mathbf{g}_i

$$\mathbf{g}_i = \sum_{n=1}^{N_{g_i}} \alpha_n^s [\mathbf{h}_n; \mathbf{t}_n],$$

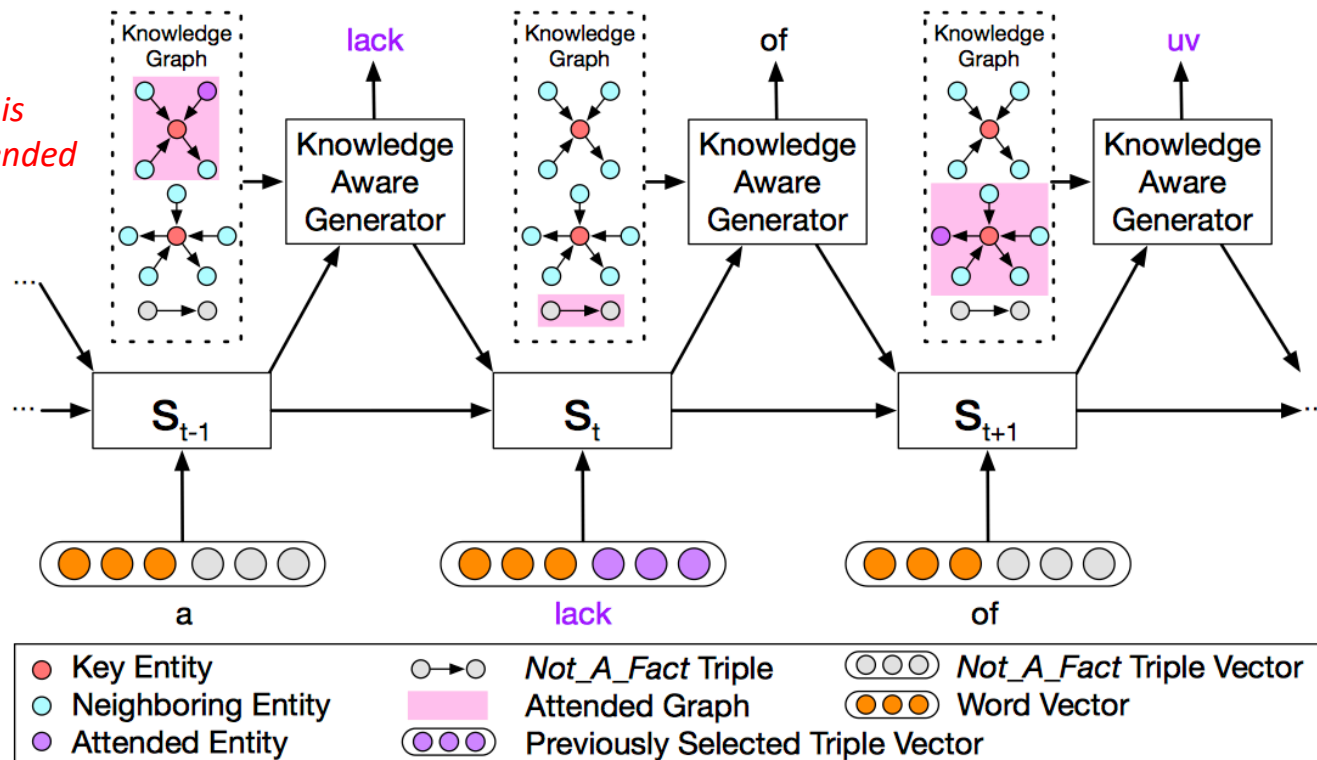
$$\alpha_n^s = \frac{\exp(\beta_n^s)}{\sum_{j=1}^{N_{g_i}} \exp(\beta_j^s)},$$

$$\beta_n^s = (\mathbf{W}_r \mathbf{r}_n)^\top \tanh(\mathbf{W}_h \mathbf{h}_n + \mathbf{W}_t \mathbf{t}_n),$$

Knowledge Aware Generator

triple's neighboring entity is used for word generation.

pink graph is mostly attended



$$s_{t+1} = \text{GRU}(s_t, [c_t; c_t^g; c_t^k; e(y_t)])$$

vectors attended on encoder hidden
vectors attended on knowledge graph vectors
vectors attended on knowledge triple vectors

$$e(y_t) = [w(y_t); k_j], \longrightarrow \text{concatenation of the word vector and the previous knowledge triple vector}$$

Dynamic Graph Attention

- Given the **decoder state** s_t , first attends on the **knowledge graph vectors** $\{g_1, g_2, \dots, g_{N_G}\}$ to compute the probability of using of each graph g_i :

the probability of choosing knowledge graph g_i at step t

$$c_t^g = \sum_{i=1}^{N_G} \alpha_{ti}^g g_i,$$

attentively reads all the knowledge graphs

$$\alpha_{ti}^g = \frac{\exp(\beta_{ti}^g)}{\sum_{j=1}^{N_G} \exp(\beta_{tj}^g)},$$

$$\beta_{ti}^g = \mathbf{V}_b^\top \tanh(\mathbf{W}_b s_t + \mathbf{U}_b g_i),$$

- The model then attends on the knowledge triple vectors $\mathbf{K}(g_i) = \{k_1, k_2, \dots, k_{N_{g_i}}\}$ within each graph g_i

$$c_t^k = \sum_{i=1}^{N_G} \sum_{j=1}^{N_{g_i}} \alpha_{ti}^g \alpha_{tj}^k k_j,$$

attentively reads all the triples in each graph

$$\alpha_{tj}^k = \frac{\exp(\beta_{tj}^k)}{\sum_{n=1}^{N_{g_i}} \exp(\beta_{tn}^k)},$$

$$\beta_{tj}^k = k_j^\top \mathbf{W}_c s_t,$$

Dynamic Graph Attention

$$\mathbf{a}_t = [\mathbf{s}_t; \mathbf{c}_t; \mathbf{c}_t^g; \mathbf{c}_t^k],$$

$$\gamma_t = \text{sigmoid}(\mathbf{V}_o^\top \mathbf{a}_t),$$

balance the choice between an entity word and a generic word

$$P_c(y_t = w_c) = \text{softmax}(\mathbf{W}_o \mathbf{a}_t),$$

P_c/P_e is the distribution over generic/entity words respectively

$$P_e(y_t = w_e) = \alpha_{ti}^g \alpha_{tj}^k,$$

$$y_t \sim \mathbf{o}_t = P(y_t) = \begin{bmatrix} (1 - \gamma_t) P_g(y_t = w_c) \\ \gamma_t P_e(y_t = w_e) \end{bmatrix},$$

Loss Function

supervised signals on the knowledge aware generator layer to teacher-force the selection of an entity or a generic word.

$$L(\theta) = - \sum_{t=1}^m \mathbf{p}_t \log(\mathbf{o}_t) - \sum_{t=1}^m (q_t \log(\gamma_t) + (1 - q_t) \log(1 - \gamma_t)),$$

$q_t \in \{0, 1\}$ *is the true choice of an entity word or a generic word in Y*

Experiments

- **Commonsense Knowledge Base: ConceptNet**
 - **world facts** such as Paris is the capital of France
 - **informal relations** between common concepts such as “A dog is a pet”
- **Commonsense Conversation Dataset:**
 - reddit single-round dialogs
 - If a post-response pair can not be connected by any triple the pair will be removed.
 - **Four test sets**
 - **High-frequency pairs:** each post has all top 25% frequent words
 - **Medium-frequency pairs:** each post contains at least one word whose frequency is within the range of 25%-75%,
 - **Low-frequency pairs:** within the range of 75%-100%
 - **OOV pairs:** each post contains out-of-vocabulary words

Baselines

- **A seq2seq model (Seq2Seq)**, which is widely used in open-domain conversational systems.
- **A knowledge-grounded model (MemNet)**, where the memory units store TransE embeddings of knowledge triples.
- **A copy network (CopyNet) model**, which copies a word from knowledge triples or generates a word from the vocabulary.

Automatic Evaluation

- **Perplexity**

- whether the content is **grammatical** and **relevant** in topic

- **Entity score**

- the number of entities per response to measure the model's ability to select the concepts from the commonsense knowledge base in generation

Model	Overall		High Freq.		Medium Freq.		Low Freq.		OOV	
	ppx.	ent.	ppx.	ent.	ppx.	ent.	ppx.	ent.	ppx.	ent.
Seq2Seq	47.02	0.717	42.41	0.713	47.25	0.740	48.61	0.721	49.96	0.669
MemNet	46.85	0.761	41.93	0.764	47.32	0.788	48.86	0.760	49.52	0.706
CopyNet	40.27	0.96	36.26	0.91	40.99	0.97	42.09	0.96	42.24	0.96
CCM	39.18	1.180	35.36	1.156	39.64	1.191	40.67	1.196	40.87	1.162

1. lowest perplexity on all the test sets
2. selects the most entities from the commonsense knowledge

rare concepts need more background knowledge to understand and respond

Manual Evaluation

- **Appropriateness**

- whether the response is appropriate in grammar, topic, and logic

- **Informativeness**

- whether the response provides new information and knowledge in addition to the post

Model	Overall		High Freq.		Medium Freq.		Low Freq.		OOV	
	app.	inf.	app.	inf.	app.	inf.	app.	inf.	app.	inf.
CCM vs. Seq2Seq	0.616	0.662	0.605	0.656	0.549	0.624	0.636	0.650	0.673	0.716
CCM vs. MemNet	0.602	0.647	0.593	0.656	0.566	0.640	0.622	0.635	0.626	0.657
CCM vs. CopyNet	0.600	0.640	0.606	0.669	0.586	0.619	0.610	0.633	0.596	0.640

Conclusion

- This work is the first attempt that uses **large-scale commonsense knowledge** in neural conversation generation.
- Instead of treating knowledge triples (or entities) separately and independently, we devise **static and dynamic graph attention mechanisms to treat the knowledge triples as a graph**

Thanks!