

The background of the slide features a silhouette of a person standing on a dark, rocky hillside. Above them is a vast, colorful sky filled with stars. The colors transition from deep purple at the top to bright yellow and orange near the horizon, creating a vibrant, aurora-like effect.

# A Hybrid Retrieval-Generation Neural Conversation Model

Hyungrak Kim

# Direction

---

- 1. Abstract & Introduction**
- 2. Related work**
- 3. Our Approach**
- 4. Experiments**
- 5. Conclusions and Future work**

# 1. Abstract & Introduction

---

- 현재 인공지능 개인 맞춤 시스템이 많은 인기를 얻고 있다
  - Amazon Alexa, Apple Siri, Microsoft Cortana
- 최근에 이런 챗봇 모델은 수작업을 줄이기 위해서 End to End Conversation 모델을 제시하고 있고 아래 대표적인 2개의 모델이 있다
- 1) Retrieval based, 2) Generation based 로 나누어 질 수 있음

# 1. Abstract & Introduction

---

- 1) Retrieval based
  - 1) 질문이 들어오면 가장 유사한 답변을 Response Repository에서 찾은 다음  
2) 가장 유사한 질문을 re-rank하고 neural ranking 모델이 최적의 답을 찾는다
  - 장점:
    - 사람의 답변에 유사한 답을 하며 Response에 대한 Controllable과 Explainable하다
    - Retrieval based 는 Response repository에 해당 질문과 가장 유사한 답변을 하기 때문에 다양한 주제에 대해서 Fluent 와 Informative 한 response를 함
  - 단점:
    - Response repository의 사이즈의 한계
    - Long tail context에 대해서 답변하기 어려움
    - Flexibility가 떨어짐 왜냐면 Response repository가 이미 Constructed 되어 있기 때문에

# 1. Abstract & Introduction

---

- 2) Generation based
  - Seq2seq모델을 바탕으로 Encoder에서 conversation context의 representation을 학습하고, Decoder에서 response sequence를 생성한다
  - 장점:
    - 어떤 Topic에도 High Coherent new Response함
  - 단점:
    - Grounding Knowledge 부족으로 인한 Generic(일반적이고) Not Informative함
      - Ex)
        - "I Have no Idea"
    - Grammar Error 존재

# 1. Abstract & Introduction

---

- 두 메소드 차이 정리

**Table 1: A comparison of retrieval-based methods and generation-based methods for data driven conversation models.**

| Item             | Retrieval-based methods   | Generation-based methods                          |
|------------------|---|---|
| Main techniques  | Retrieval models; Neural ranking models                                   | Seq2Seq models                                    |
| Diversity        | Usually good if similar contexts have diverse responses in the repository | Easy to generate bland or universal responses     |
| Response length  | Can be very long  | Usually short                                     |
| Context property | Easy for similar context in the repository; Hard for unseen context       | Easy to generalize to unseen context              |
| Efficiency       | Building index takes long time; Retrieval is fast                         | Training takes long time; Decoding is fast        |
| Flexibility      | Fixed response set once the repository is constructed                     | Can generate new responses not covered in history |
| Fluency          | Natural human utterances  | Sometimes bad or contain grammar errors           |
| Bottleneck       | Size and coverage of the repository                                       | Specific responses; Long text; Sparse data        |
| Informativeness  | Easy to retrieve informative content                                      | Hard to integrate external factual knowledge      |
| Controllability  | Easy to control and explain   | Difficult to control the actual generated content |

# 1. Abstract & Introduction

---

- 이 논문에서는 두가지 기술의 Merits를 결합한 Hybrid neural conversation model을 제시
  - 1) Generation model에서 Response를 생성하고(Seq2seq모델),
  - 2) Retrieval Model에서 Response candidates의 set을 Recall하고,
  - 3) neural ranking model이 1)과 2)으로부터 best response candidate를 select한다(neural ranking model은 conversation pair의 representation과 matching feature를 학습)
    - 이를 학습하기 위해서 distante supervision을 제안하는데, retrieval 모델/generation모델에서 생성한 Response Candidates의 (Positive, Negative)label을 자동으로 주론 해서 Neural Ranking Model의 학습데이터를 제공

# 1. Abstract & Introduction

---

- 실험 결과는 Twitter와 Foursquare data에서 다른 2가지 모델보다 Outperform 하는 것을 보여줌
- Contribution
  - Retrieval-based model과 Generation-based model을 비교한 것
  - Hybrid neural conversational 모델을 제안한 것
  - Distant Supervision approach 제시
  - 자동 평가와 사람 평가에서 좋은 결과를 가져옴

## 2. Related work

---

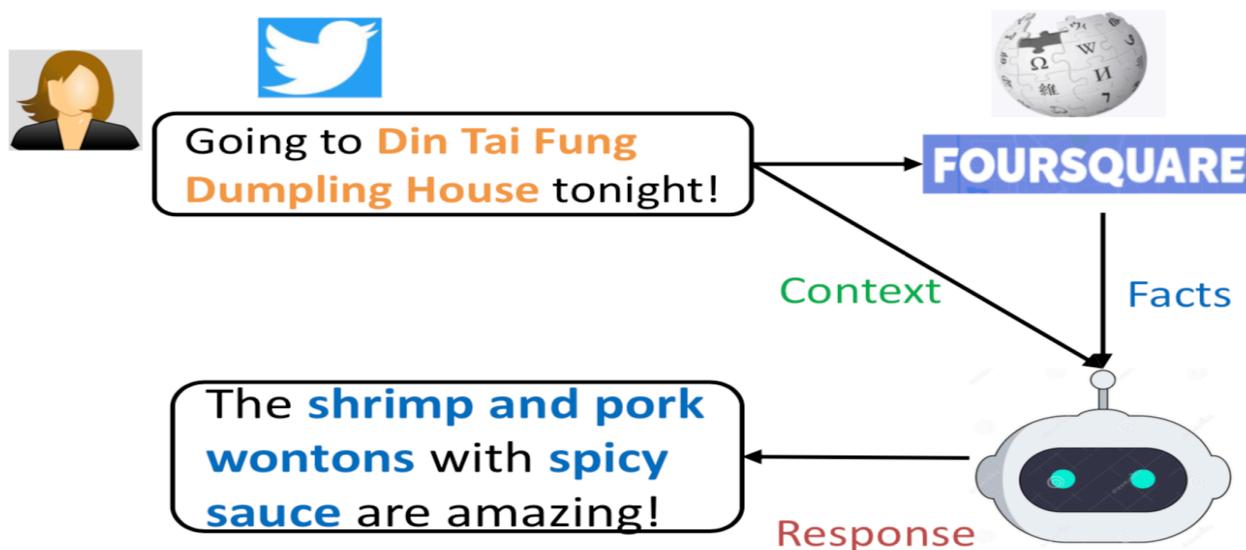
- Neural Ranking Models
  - 3 Category
    - 1. Representation-focused models
      - Queries와 Documents의 Representation를 학습하고 이 둘의 유사도를 계산
    - 2. Interaction-focused model
      - Query-document의 쌍의 interaction matrix 계산해서 매칭
    - 3. 1, 2번 neural ranking 모델을 혼합해서 사용 Lexical 과 semantic matching을 사용
  - 이 Approach에 사용한 모델은 Interaction-focused 모델에 속함

# 3. Our Approach

---

- 3.1 Problem Formulation

- Conversation Context에서  $u_i$  is the i-th context sequence
- F factual snippets of text  $F_i = \{f_{i1}, f_{i2}, \dots, f_{iF}\}$



**Figure 1: An example of the conversational response generation task. The factual information from external knowledge is denoted in blue color.**

# 3. Our Approach

---

- 3.2 Method Overview
  - Hybrid Neural Conversation Model(HybridNCM)
    - 1)Generation module
      - Conversation Context  $u_i$  와  $F_i$  (Fact)를 Seq2seq 모델에서 Response Candidate  $G_i$  Set을 생성
      - Context Encoder
        - Seq2seq model + attention mechanism
        - $h_t = RNN(u_t, h_{t-1})$ ,
        - 2 LSTM Layer
      - Fact Encoder
        - Context Encoder와 같은 Layer
      - Response Decoder
        - Next Word  $G_t$ 를 생성, 2 LSTM Layer + attention

# 3. Our Approach

---

- 4) -> 이전 Decoder hidden state St-1 의 Attention at
- 8) Encoder의 Last output 을 가지고 Decoder 초기화

$$p(g|u_i, \mathcal{F}) = \prod_{t=1}^{L_g} p(g_t | g_{1:t-1}, u_i, \mathcal{F}) \quad (2)$$

$$\mathbf{E} = [\mathbf{h}_1, \dots, \mathbf{h}_{L_u}, \bar{\mathbf{f}}^1, \dots, \bar{\mathbf{f}}^F] \in \mathbb{R}^{H \times (L_u + F)} \quad (3)$$

$$\mathbf{a}_t = \text{softmax}(\mathbf{E}^T \mathbf{s}_{t-1}) \quad (4)$$

$$\mathbf{c}_t = \mathbf{E} \mathbf{a}_t \quad (5)$$

$$\mathbf{v}_t = \tanh([\mathbf{s}_{t-1}, \mathbf{c}_t]) \quad (6)$$

$$\mathbf{s}_t = \text{RNN}(\mathbf{v}_t, \mathbf{s}_{t-1}) \quad (7)$$

$$\mathbf{s}_0 = \varphi \left( \tanh \left( \mathbf{h}_{L_u} + \frac{1}{F} \sum_{j=1}^F \bar{\mathbf{f}}^j \right) \right) \quad (8)$$

- 이전 단어 다음에 나올 단어 예측  $p(g_t | g_{1:t-1}, u_i, \mathcal{F}) = \text{softmax}(\phi([\mathbf{s}_{t-1}, \mathbf{c}_t]))$

# 3. Our Approach

---

- Train and Decode

$$\mathcal{L}_g = -\frac{1}{|\mathcal{U}|} \sum_{y^*, u_i, \mathcal{F}} \log p(y^* | u_i, \mathcal{F}) \quad (10)$$

- Negative log-likelihood 를 모든 트레이닝 데이터에 대해서 훈련
- Response를 생성하기 위해서 Beam search를 이용
- 짧은 답변을 생성한 경우에 Penalty를 주기위해 log-likelihood 점수를 생성된 단어를 length로 나누어 준다

# 3. Our Approach

---

- Hybrid Neural Conversation Model(HybridNCM)
  - 2) Retrieval Module
    - “Context-Context Match” Approach
    - 모든 historical conversation context에서 현재 conversation context  $u_i$ 를 매칭하여 Response candidate  $R_i$ 를 생성

# 3. Our Approach

---

- Hybrid Neural Conversation Model(HybridNCM)
  - 3) Hybrid Ranking Module
    - Interaction Matching Matrix
      - $Y_i = G_i \cup R_i$
      - Generated Response  $G_i$  Set
      - Response Candidate  $R_i$  Set
      - Hybrid Ranking Candidate Re-Rank all candidate  $Y_i$
      - Conversation Context  $u_i$
    - Distant Supervision으로부터 라벨로 학습된 Hybrid Neural Ranker 를 통해서 all candidate  $Y_i$  Best response를 출력함
    - Fact는 사용하지 않고
    - Conversation Context  $u_i$  all candidate  $Y_i$  의 Interaction Matching MatrixCNN과 MLP로 계산

# 3. Our Approach

---

- Hybrid Neural Conversation Model(HybridNCM)
  - Ranking을 매기 위한 데이터 구성
    - triples( $u_i, y_{k+}, y_{k-}$ ), where  $y_{k+}$  and  $y_{k-}$  denote the positive and the negative response candidate 를 학습 데이터로 사용
  - 입력 Context와 적절한 Response를 평가하는 것이 쉽지 않고, 또 비용이 많이 들어가므로 입력된 Context에 대한 적절한 응답을 계산하는 것이 필요
  - 따라서 효과적으로 트레이닝 하기 위해서 Ranking module을 학습하기 위한 Training data를 생성할 필요가 있음
  - Hybrid ranking model에서 생성한 Response Candidate  $y_i$  와 BLEU/ROUGE-L 과 같은 기법으로 Ground Truth를 비교한다. 이 비교를 통해서 Positive Candidate response와 정답과 다른 Negative candidate를 구별하고 Training data의 라벨로 만듬
  - Hybrid Ranking model의 Hinge Loss

$$\mathcal{L}_h = \sum_{i=1}^I \max(0, \epsilon - f(u_i, y_i^{k+}) + f(u_i, y_i^{k-})) + \lambda \|\Theta\|_2^2 \quad (13)$$

- Relation Extraction 기법에 영감을 받음
- Negative 개념이 들어가는건 더 좋은 Response를 구별하기 위해서

# 3. Our Approach

- Total Architecture

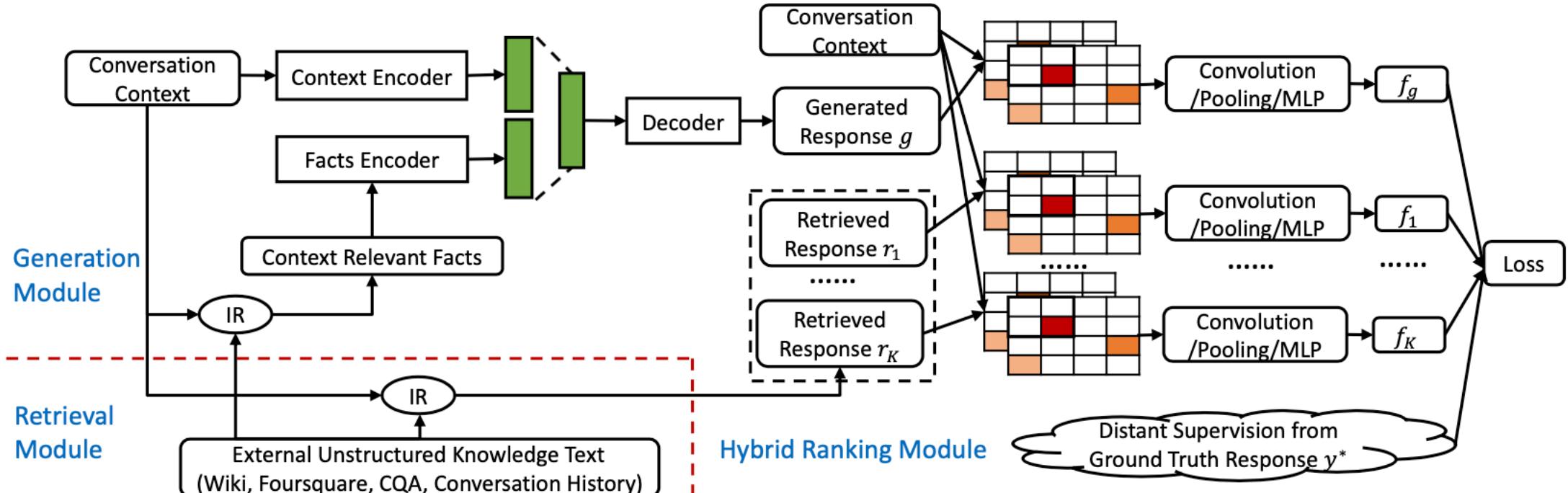


Figure 2: The architecture of the Hybrid Neural Conversation Model (HybridNCM).

# 4. Experiments

---

- 4.1 Data Set Description
  - Twitter, Foursquare Dataset

**Table 3: Statistics of experimental data used in this paper.**

| Items                    | Train      | Valid  | Test   |
|--------------------------|------------|--------|--------|
| # Context-response pairs | 1,059,370  | 2,067  | 2,066  |
| # Facts                  | 43,111,643 | 79,950 | 79,915 |
| Avg # facts per context  | 40.70      | 38.68  | 38.68  |
| Avg # words per facts    | 17.58      | 17.42  | 17.47  |
| Avg # words per context  | 16.66      | 17.85  | 17.66  |
| Avg # words per response | 11.65      | 15.58  | 15.89  |

# 4. Experiments

---

- 4.2 Experimental Setup
  - 3가지 모델을 비교
    - Retrieval-based
    - Generation-based
    - Hybrid retrieval-generation method
      - 1) HybridNCM-RS
      - 2) HybridNCM-RSF(Fact Encoder 추가)
  - Seq2seq
  - Seq2seq-Fact(fact encoder 추가),
  - KNCM-MTask-R(Knowledge-grounded neural conversation model, 23 million general Twitter conversation, 1 million fact data) 최근 가장 좋은 모델
  - Evaluation
    - 생성된 response에 대한 BLEU와 ROUGE-L기법을 사용
    - Corpus Level에서 BLEU평가가 Sentence Level평가보다 더 사람의 평가에 상관성이 있다.
    - Lexical 에 diversity(Distinct-1과 Distinct-2 사용) 와 informativeness 를 평가
    - 사람의 평가도 들어감

# 4. Experiments

---

- 4.2 Experimental Setup
  - Parameter Settings

| Models                           | Seq2Seq | Seq2Seq-Facts |
|----------------------------------|---------|---------------|
| Embedding size                   | 512     | 256           |
| # LSTM layers in encoder/decoder | 2       | 2             |
| LSTM hidden state size           | 512     | 256           |
| Learning rate                    | 0.0001  | 0.001         |
| Learning rate decay              | 0.5     | 0.5           |
| # Steps between validation       | 10000   | 5000          |
| Patience of early stopping       | 10      | 10            |
| Dropout                          | 0.3     | 0.3           |

# 4. Experiments

---

- 4.3 Evaluation Results

| Method        | BLEU          | ROUGE-L                    | Distinct-1    | Distinct-2    |
|---------------|---------------|----------------------------|---------------|---------------|
| Seq2Seq       | 0.5032        | 8.4432                     | 2.36%         | 11.18%        |
| Seq2Seq-Facts | 0.5904        | 8.8291                     | 1.91%         | 7.85%         |
| KNCM-MTask-R  | 1.0800        | \                          | 7.08%         | 21.90%        |
| Retrieval     | 1.2491        | 8.6302                     | <b>14.68%</b> | <b>58.71%</b> |
| HybridNCM-RS  | 1.3450        | <b>10.4078<sup>‡</sup></b> | 11.30%        | 47.35%        |
| HybridNCM-RSF | <b>1.3695</b> | 10.3445 <sup>‡</sup>       | 11.10%        | 46.01%        |

# 4. Experiments

---

- 4.3 Evaluation Results

- 사람 평가

- Appropriateness
    - Informativeness
    - 0(bad), +1(neutral), +2(good)

| Comparision   | Appropriateness           |        |        |             |          | Informativeness           |        |        |             |          |
|---------------|---------------------------|--------|--------|-------------|----------|---------------------------|--------|--------|-------------|----------|
|               | Method                    | Mean   | Bad(0) | Neutral(+1) | Good(+2) | Agreement                 | Mean   | Bad(0) | Neutral(+1) | Good(+2) |
| Seq2Seq       | 0.4733                    | 61.67% | 29.33% | 9.00%       | 0.2852   | 0.2417                    | 77.58% | 20.67% | 1.75%       | 0.4731   |
| Seq2Seq-Facts | 0.4758                    | 62.50% | 27.42% | 10.08%      | 0.3057   | 0.3142                    | 70.75% | 27.08% | 2.17%       | 0.4946   |
| Retrieval     | 0.9425                    | 34.42% | 36.92% | 28.67%      | 0.2664   | 0.8008                    | 35.50% | 48.92% | 15.58%      | 0.3196   |
| HybridNCM-RS  | <b>1.1175<sup>‡</sup></b> | 27.83% | 32.58% | 39.58%      | 0.3010   | <b>1.0650<sup>‡</sup></b> | 18.42% | 56.67% | 24.92%      | 0.1911   |
| HybridNCM-RSF | 1.0358 <sup>‡</sup>       | 31.67% | 33.08% | 35.25%      | 0.2909   | 1.0292 <sup>‡</sup>       | 20.42% | 56.25% | 23.33%      | 0.2248   |

# 4. Experiments

---

- 4.3 Evaluation Results

- 사람 평가

| Type                    | Appropriateness | Informativeness |
|-------------------------|-----------------|-----------------|
| Comparision             | Win/Tie/Loss    | Win/Tie/Loss    |
| HNCM-RS v.s. Seq2Seq    | 0.71/0.15/0.14  | 0.84/0.10/0.06  |
| HNCM-RSF v.s. Seq2Seq   | 0.68/0.16/0.16  | 0.82/0.11/0.07  |
| HNCM-RS v.s. Seq2Seq-F  | 0.70/0.15/0.15  | 0.80/0.12/0.08  |
| HNCM-RSF v.s. Seq2Seq-F | 0.65/0.19/0.17  | 0.77/0.15/0.09  |
| HNCM-RS v.s. Retrieval  | 0.43/0.31/0.26  | 0.50/0.31/0.18  |
| HNCM-RSF v.s. Retrieval | 0.41/0.30/0.29  | 0.50/0.28/0.22  |

# 4. Experiments

---

- 4.4 Analysis of Top Responses Selected By Re-Ranker

| Item            | HybridNCM-RS | HybridNCM-RSF |
|-----------------|--------------|---------------|
| #TestQNum       | 2066         | 100.00%       |
| #PickedGenRes   | 179          | 8.66%         |
| #PickedRetRes   | 1887         | 91.34%        |
| #PickedTop1BM25 | 279          | 13.50%        |
|                 | 2066         | 100.00%       |
|                 | 275          | 13.31%        |
|                 | 1791         | 86.69%        |
|                 | 253          | 12.25%        |

# 4. Experiments

---

- 4.5 Impact of Distant Supervision Signals

| Model       | HybridNCM-RS  |                | HybridNCM-RSF |                |
|-------------|---------------|----------------|---------------|----------------|
| Supervision | BLEU          | ROUGE-L        | BLEU          | ROUGE-L        |
| BLEU-1      | <b>1.3450</b> | <b>10.4078</b> | <b>1.3695</b> | <b>10.3445</b> |
| BLEU-2      | 1.1165        | 10.1584        | 0.8239        | 9.8575         |
| ROUGE-L     | 1.1435        | 10.0928        | 0.9838        | 9.7961         |
| SentBLEU    | 0.8326        | 9.2887         | 1.0631        | 9.6338         |

# 4. Experiments

---

- 4.6 Impact of Ratio of Positive Samples

|               | Supervision | BLEU-1        |                | BLEU-2        |                | ROUGE-L       |                |
|---------------|-------------|---------------|----------------|---------------|----------------|---------------|----------------|
| Model         | # Positive  | BLEU          | ROUGE-L        | BLEU          | ROUGE-L        | BLEU          | ROUGE-L        |
| HybridNCM-RS  | k'=1        | 0.9022        | 8.9596         | 0.7547        | 8.8351         | 1.0964        | 8.9234         |
|               | k'=2        | 1.0649        | 9.7241         | 1.1099        | 9.9168         | 1.1019        | 9.6216         |
|               | k'=3        | <b>1.3450</b> | <b>10.4078</b> | <b>1.1165</b> | <b>10.1584</b> | <b>1.1435</b> | <b>10.0928</b> |
| HybridNCM-RSF | k'=1        | 1.0223        | 9.2996         | <b>1.1027</b> | 9.2453         | 1.0035        | 9.2812         |
|               | k'=2        | 1.3284        | 9.8637         | 1.0175        | 9.8562         | <b>1.0999</b> | <b>9.8061</b>  |
|               | k'=3        | <b>1.3695</b> | <b>10.3445</b> | 0.8239        | <b>9.8575</b>  | 0.9838        | 9.7961         |

# 4. Experiments

---

- 4.7 Example and Case Study

| Context       | Donated to the [ENTITY] last night and now I have to listen to automated phone calls. It's enough to make me want to cancel. |  |
|---------------|--|--|
| Method        | r/g  | System Output Response   |
| Ground Truth  | -  | Ask them to put you on their internal dnc list. They will likely respect this, because future calls can get them charged.              |
| Seq2Seq       | g  | I didn't get it. I didn't.   |
| Seq2Seq-Facts | g  | I'm sorry to hear that. Please dm us your email address so we can look into this. Thanks!  |
| Retrieval     | r  | It's a known issue in LA county. I just got an email from my dm and tech and it should be good in 30 minutes or so.                    |
| HNCM-RS       | r  | We're listening and would like to know more and help with your experience. Please follow us so i can dm you our contact info. [ENTITY] |
| HNCM-RSF      | g  | We're sorry to hear this. Please dm us if you need assistance. Please dm us your contact info so we can look into this.                |

# 5. Conclusions and Future work

---

- Retrieval 과 Generation을 혼합하여 사용하고, 학습을 할때 Distant Supervision Approach를 이용해서 학습데이터를 자동으로 label화를 시킨다.
- 이를 통해서 실험에서 위와 같은 Approach가 효과가 있다고 증명하고 있다

---

# Q&A

---

**Thank you**