



Incorporating Behavioral Hypotheses for Query Generation

Ruey-Cheng Chen (SEEK Ltd)
Chia-Jung Lee (Microsoft)

EMNLP 2020 / Short Paper

[paper](#) / [talk](#)

집현전 중급반
발표자 송일현

Incorporating Behavioral Hypotheses for Query Generation



≈ Query Suggestion, Related Query

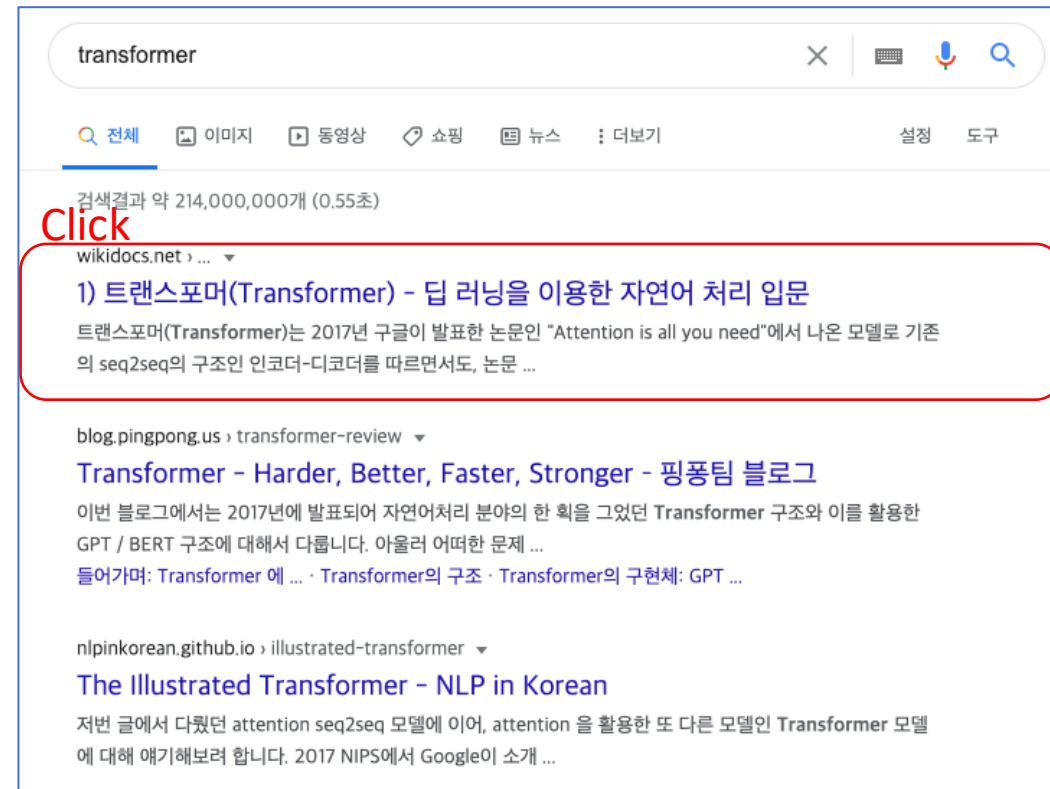
The screenshot displays a job search interface with three main components illustrating query generation:

- Left Panel (seek Job Search):** A search bar contains "machine learning". Below it, a dropdown menu lists suggestions: "machine learning", "machine learning engineer", "machine learning developer", and "junior machine learning".
- Middle Panel (Google):** A search bar contains "machine learning". Below it, a list of suggestions is shown: "machine learning", "machine learning for kids", "machine learning algorithm", "machine learning deep le", "machine learning model", "machine learning pdf", "machine learning engineer", "machine learning a proba", and "machine learning이란".
- Right Panel (N 머신러닝):** A search bar contains "머신러닝". Below it, a list of suggestions is shown: "5기 연봉평균 4050만원. 과장광고 없이 결과로 얘기하는 코딩교육기관", "Dataiku www.hancomit.com", and "데이터머신러닝을 위한,데이터머신러닝 협업 플랫폼Dataiku". Below the suggestions, a "더보기 →" button is visible. At the bottom, a "연관 검색어" (Related Search) section lists: "머신러닝 딥러닝 차이", "딥러닝", "machine learning", "ai 기술", "인공지능 딥러닝", "ai 딥러닝", "기계학습", "인공지능 기술", "인공지능 역사", and "딥러닝 기술".

Incorporating **Behavioral Hypotheses** for Query Generation



- User Behavior
 - Query -> Click (Document : title)
 - Transformer
 - > "1) **트랜스포머**(Transformer) – **딥러닝**을 이용한 **자연어처리**..."
- Query -> Query
 - Machine Learning Engineer
 - > Machine Learning Developer

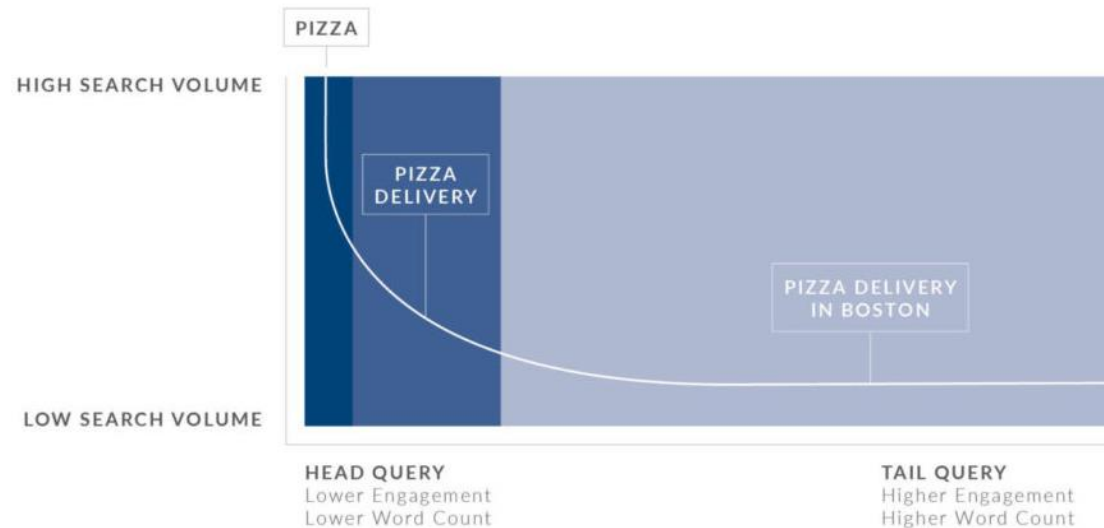


Introduction



- Query Suggestion
 - **Discriminative** Approach
 - 사용자 로그에서 다음에 발생할 확률이 높은 질의를 **선택** (Query Co-occurrence)
 - 로그에 있는 질의만 처리 가능, **tail query** 에 penalty

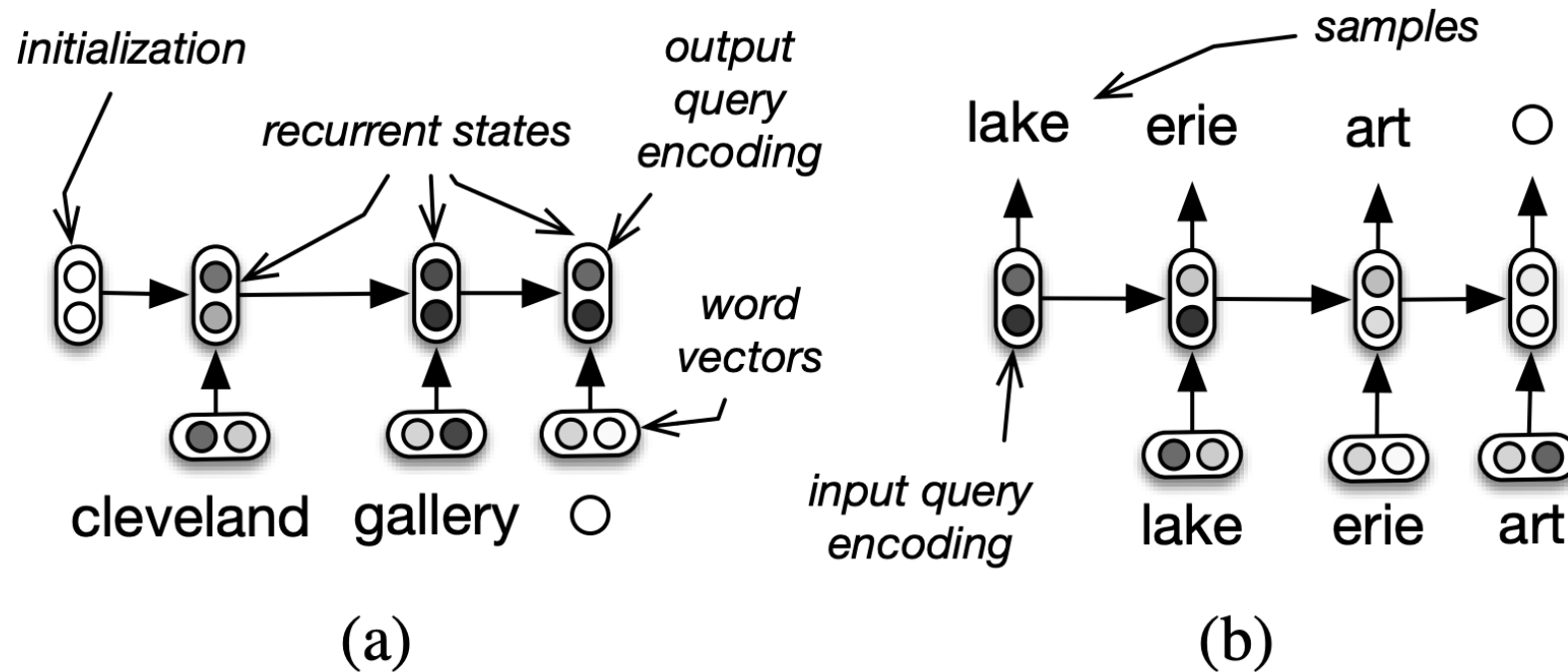
Search Volume vs. Engagement vs. Number of Words





Introduction

- Query Suggestion
 - **Generative** Approach
 - 질의를 **생성** (ex. RNN Encoder-Decoder)

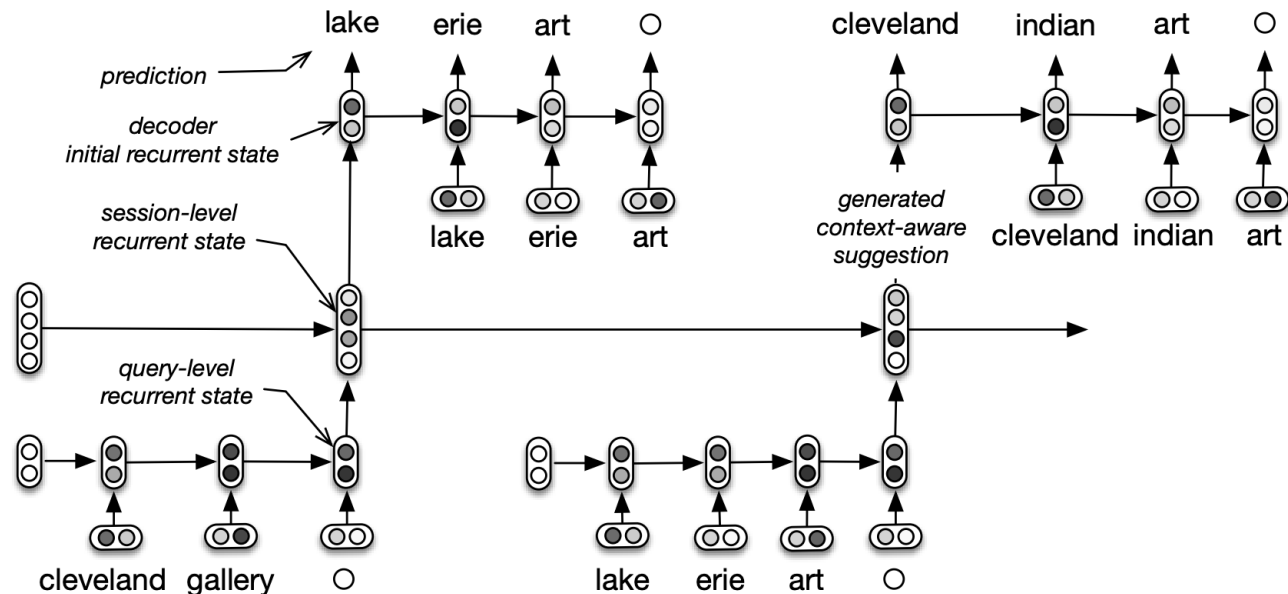


cleveland gallery -> lake erie art



Introduction

- Query Suggestion
 - Generative Approach
 - **HRED** (Hierarchical Recurrent Encoder-Decoder)



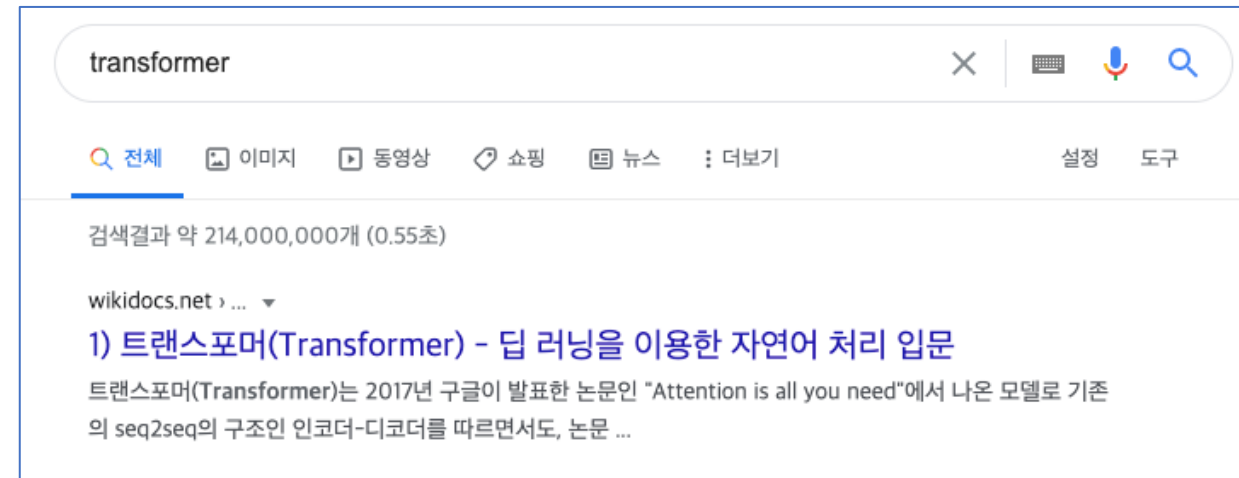
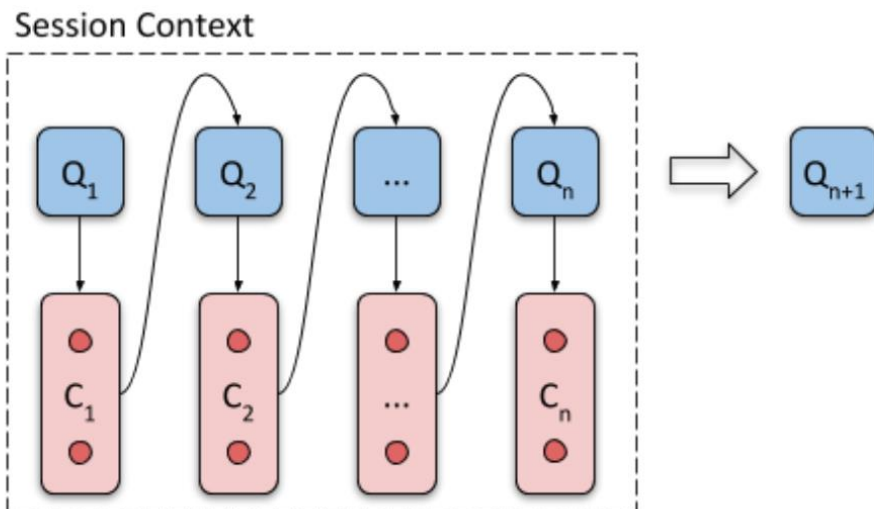
[cleveland gallery -> lake erie art] -> cleveland indian art

한계 : 질의-질의 간의 관계만 반영
=> User의 Implicit 한 Intent(click) 정보도 포함 해보자



Approach

- $Q = (Q_1, Q_2, \dots, Q_n)$: 사용자가 입력한 질의 Sequence (in 세션)
- $C = (C_1, C_2, \dots, C_n)$: $C_i \rightarrow Q_i$ 의 결과에서 사용자의 Interaction
ex) 사용자가 클릭한 문서의 제목
 - Q_i : Transformer
 - C_i : **트랜스포머**(Transformer) – **딥러닝**을 이용한 **자연어 처리** 입문
- Q, C 를 이용해서 Q_{n+1} (다음 질의) 생성





Approach : Behavioral Hypotheses

- Q_{n+1} 생성하는데 무엇이 중요할까? (가정)

K_1 : 사용자가 **입력한 질의**들이 중요할 것이다.

$$K_1 = (Q_1, Q_2, \dots, Q_n)$$

K_2 : 사용자가 **반응(Click)한 문서**들과 **직전 질의**가 중요할 것이다.

$$K_2 = \bigoplus_{i < n} \{(t)_{t \in C_i}\} \oplus (Q_n)$$

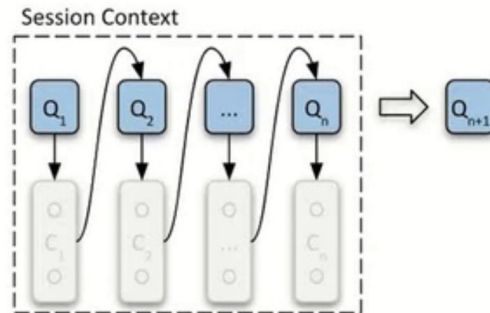
K_3 : 사용자가 **반응한 질의와 문서**들이 중요할 것이다.

$$K_3 = \bigoplus_{Q_i: C_i \neq \emptyset} \{(Q_i) \oplus (t)_{t \in C_i}\} \oplus (Q_n)$$

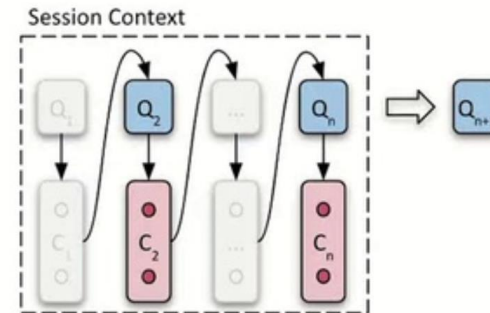
K_4 : 사용자의 **직전에 입력한 질의**와 **반응한 문서**가 중요할 것이다.

$$K_4 = (Q_n) \oplus (t)_{t \in C_n}$$

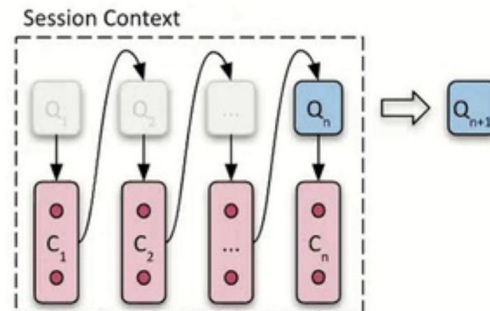
K_1 : All queries



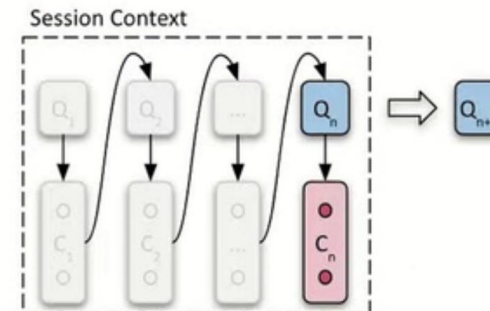
K_3 : Interacted queries/cues + last query



K_2 : All matching cues + last query



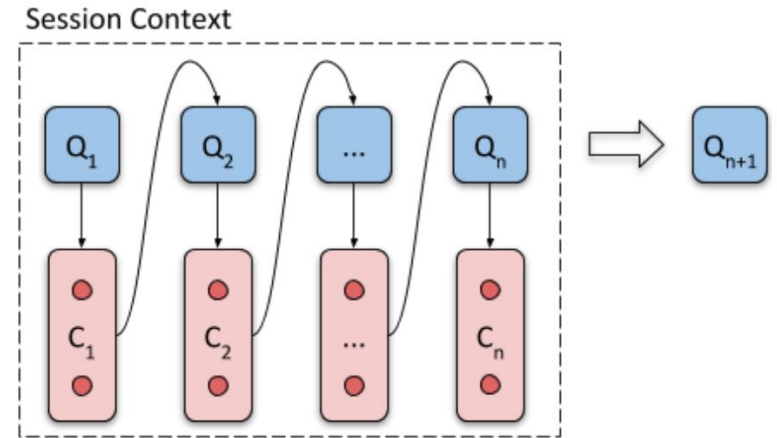
K_4 : Last query + cues



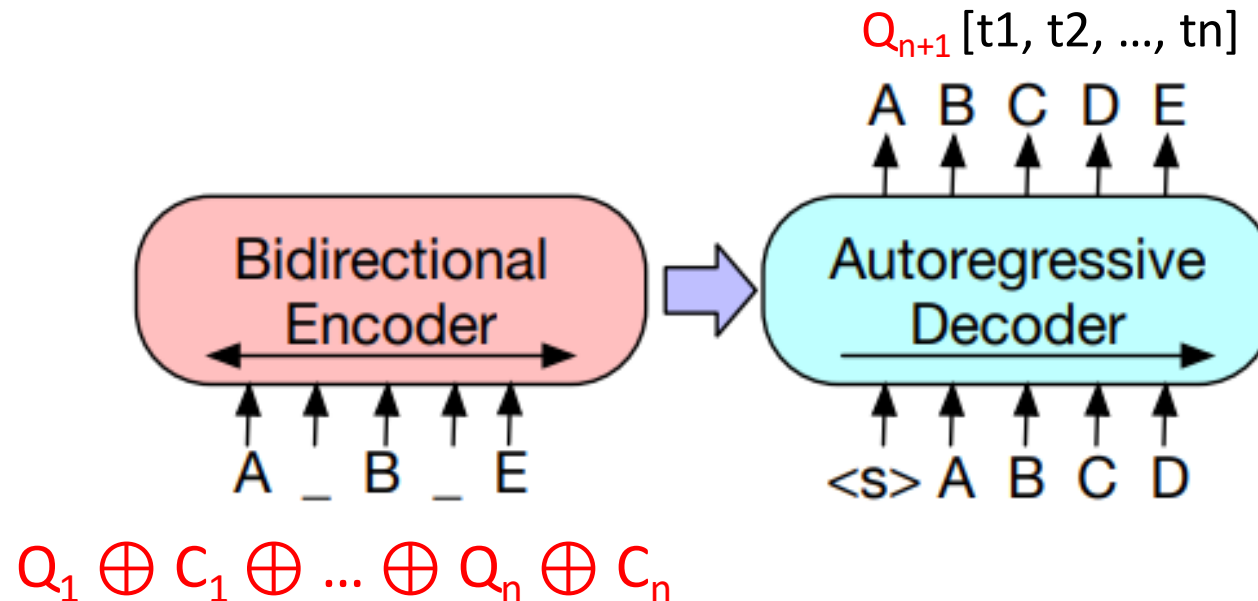
most recent interactions are the most important

Approach : Vanilla Encoder-Decoder Transformer

- 이 가정들을 어떻게 사용할까?
 - Vanilla Encoder-Decoder Transformer
 - User Behavior 를 **Concatenate** 한 입력에서 Q_{n+1} 을 생성 (번역 모델과 유사한 Seq2Seq)
 - $Q_1 \oplus C_1 \oplus Q_2 \oplus C_2 \oplus \dots \oplus Q_n \oplus C_n \rightarrow Q_{n+1}$



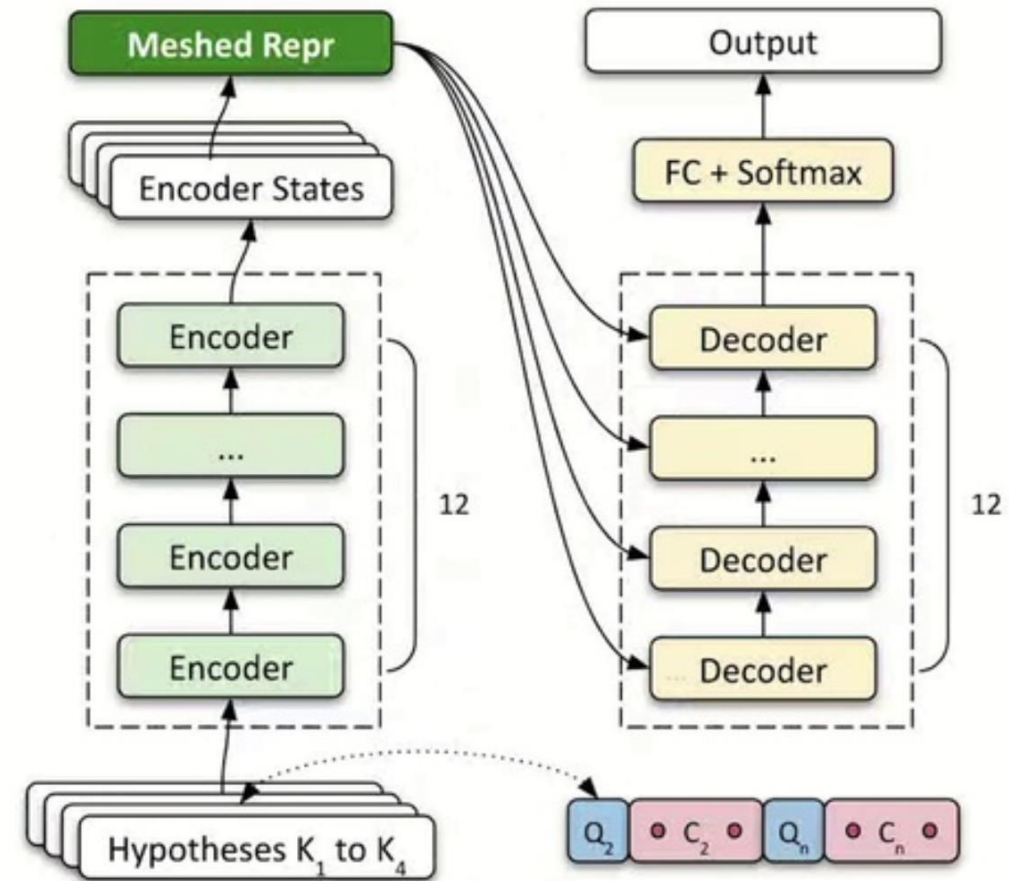
- BART** 사용



Approach : Meshed Representations



- 이 가정들을 어떻게 사용할까?
- Meshed Representations
 - $K_1 \sim K_4$ 의 encoding된 token 들을 **Token wise Attention**를 통해 Mesh
 - **Meshed Repr**의 결과를 Decoder로 전달



Token wise Attention

$$K_1 = (Q_1, Q_2, \dots, Q_n)$$

$$K_2 = \bigoplus_{i < n} \{(t)_{t \in C_i}\} \oplus (Q_n)$$

$$K_3 = \bigoplus_{Q_i: C_i \neq \emptyset} \{(Q_i) \oplus (t)_{t \in C_i}\} \oplus (Q_n)$$

$$K_4 = (Q_n) \oplus (t)_{t \in C_n}$$

The procedure can be described as follows. Let $[S_i^{(1)}; \dots; S_i^{(T)}] = \text{BART}_{\text{enc}}(K_i)$ for all K_i , and T is the sequence length. We have:

$$\alpha_i^{(j)} \propto \exp(W_{\text{attn}} S_i^{(j)})$$

$$F^{(j)} = \sum_i \alpha_i^{(j)} S_i^{(j)} \quad (2)$$

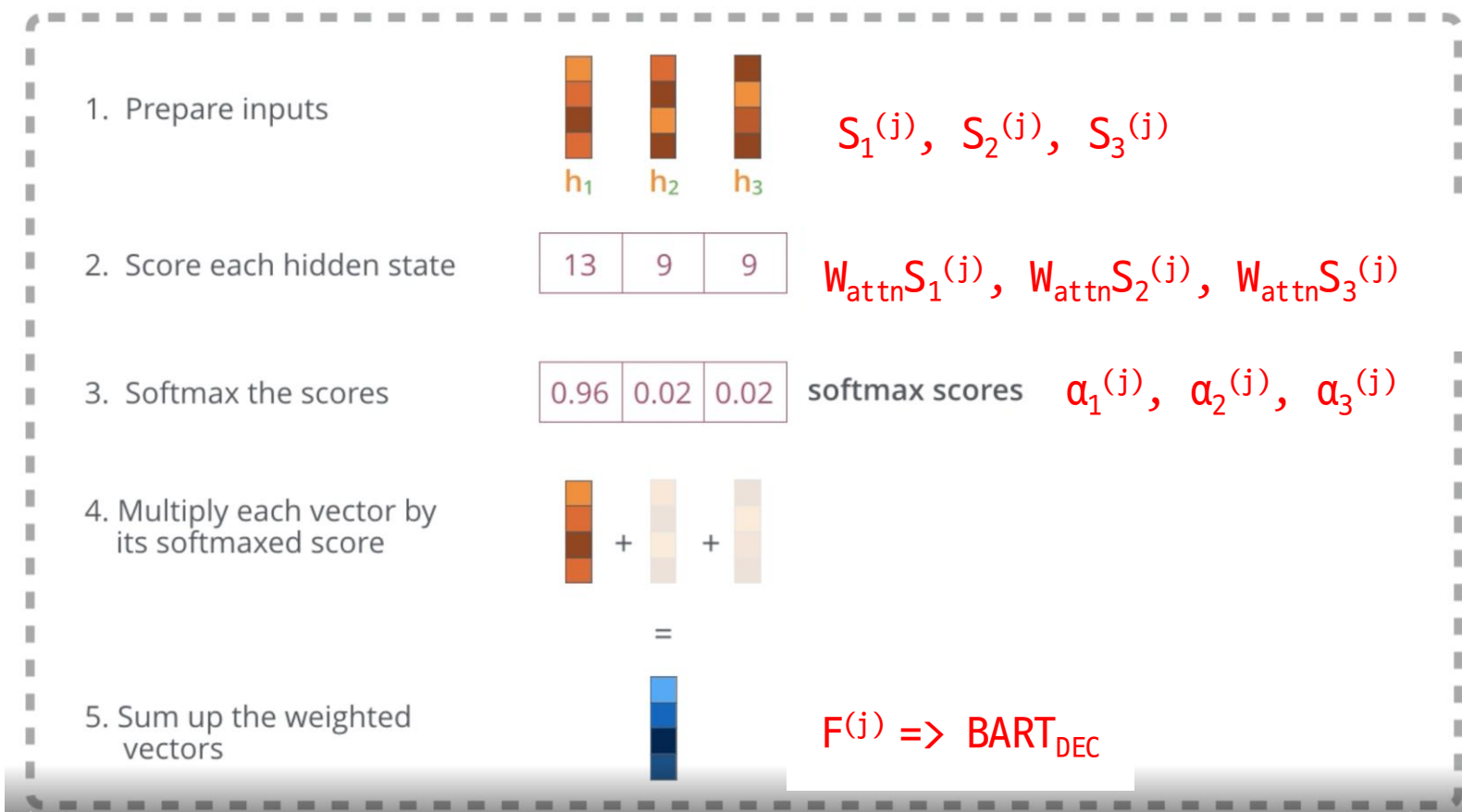
$$O = \text{BART}_{\text{dec}}([F^{(1)}; \dots; F^{(T)}])$$

where W_{attn} is the attention weight matrix to be learned and O the output. On the decoder side, the

<https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>



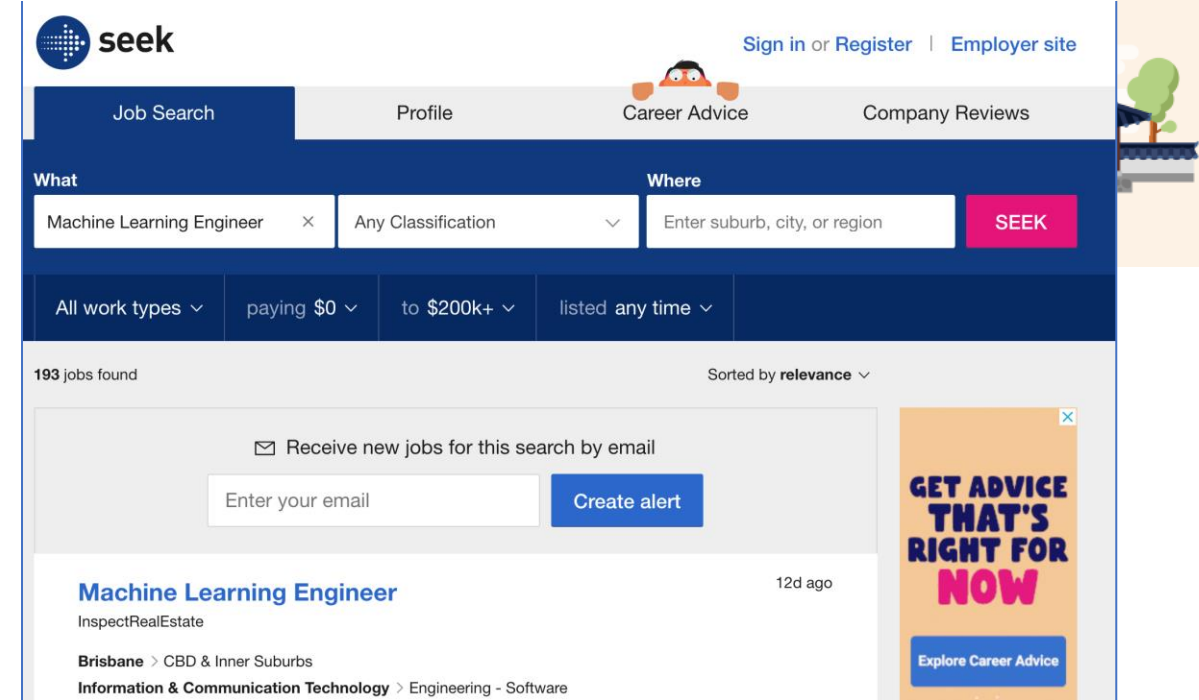
$S_1^{(j)}$ → K_1 의 j 번째 token의 BART Encoding → h_1
 $S_2^{(j)}$ → K_2 의 j 번째 token의 BART Encoding → h_2
 $S_3^{(j)}$ → K_3 의 j 번째 token의 BART Encoding → h_3
 $S_4^{(j)}$ → K_4 의 j 번째 token의 BART Encoding → h_4 (그림생략)



$\text{Dim}(F^{(j)}) == \text{Dim}(S_i^{(j)})$: Pretrained Model 사용 가능!

Experiment : Data

- [SEEK] : Job search engine in AU
 - Role Title, Skill, Company Name, Geo Location
 - Search Session log
- Q_i : Query
- C_i : **Title** of documents that were clicked on in response to Q_i
- Session boundary: inactivity of 30 minutes or more between two consecutive actions.
- In each session **the latest query** was held out as the **ground truth**.
- Training Session (500K) + DEV(1K) : 2019'10 월 2주 / Test (100K) : 그 후 2주
- Data Cleansing
 - BART의 maximum sequence length 넘어가는 15% 세션 제거
 - 10회 이하 출현한 noisy query 제거
 - Singleton session (contain only one query) 제거

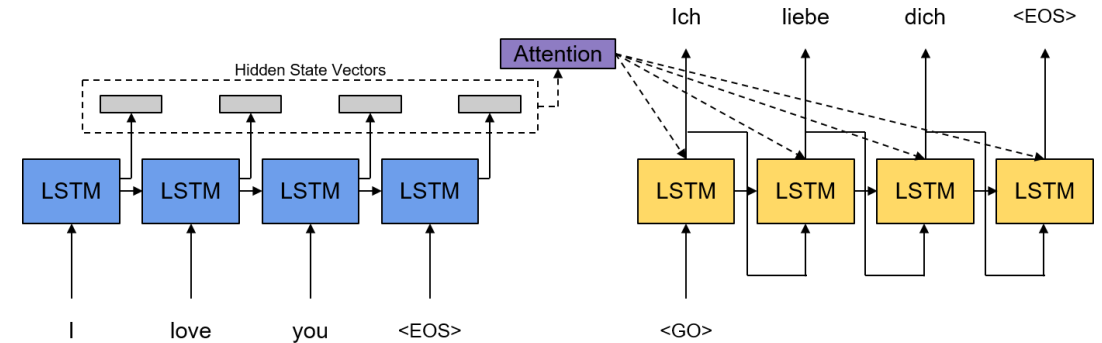


Experiment : 비교 모델



<https://hcnoh.github.io/2018-12-11-bahdanau-attention>

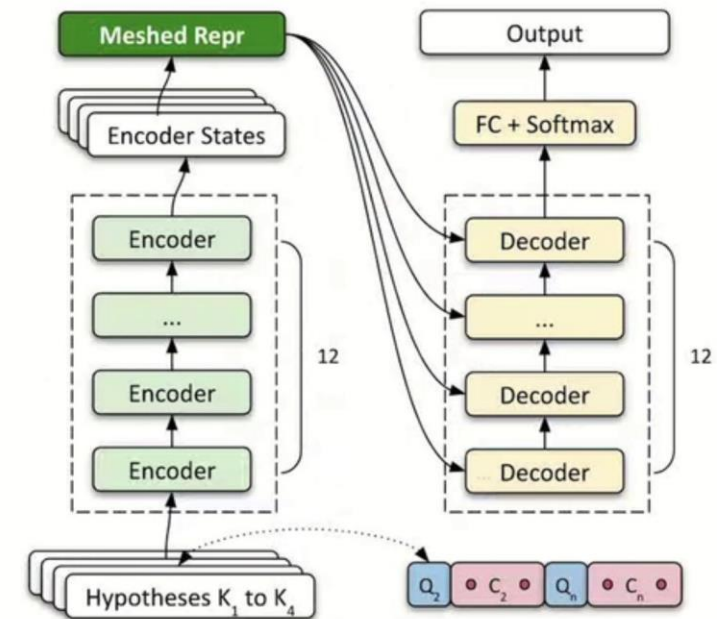
- **Seq2Seq-Attn** (Bahdanau et al., 2014)
 - GRU, 1000 hidden dim, byte-pair encoding



- **MPS (Most Popular Suggestions)**
 - **Co-occurrence Frequencies** of the last query in the search context and all candidate queries.

ML Scientist -> ML Engineer (10번),
ML Scientist -> Data Scientist (4번)

- **BART - Vanilla BART**
 - Full concatenate search context as INPUT
 - $Q_1 \oplus C_1 \oplus Q_2 \oplus C_2 \oplus \dots \oplus Q_n \oplus C_n \rightarrow Q_{n+1}$



- **Mesh BART**- 제안 모델 (p.10)

Experiment : Metric (1/2)



- **Word Error Rate** $WER@k = \min_{i=1, \dots, k} \text{EditDist}(ref, hyp^{(i)}) / |ref|$

EditDist("Lead Data Engineer", "Data Scientist") = 2

(Delete:Lead, Change:Engineer -> Scientist)

- **Mean Reciprocal Rank (MRR@K)** - https://en.wikipedia.org/wiki/Mean_reciprocal_rank
Rank 정답이 있는 순위 (1~K)

Reciprocal Rank = $1/\text{rank}$

정답이 1등에 있음 - $1/1 = 1$

2등에 있음 $1/2 = 0.5$

3등에 있음 $1/3 = 0.333$

K등안에 없음 0

- **Success at K (S@K)** : K등 안에 정답 있음

Experiment : Metric (2/2)



- **BERT F1**– BERTScore: Evaluating Text Generation with BERT / <https://arxiv.org/abs/1904.09675>

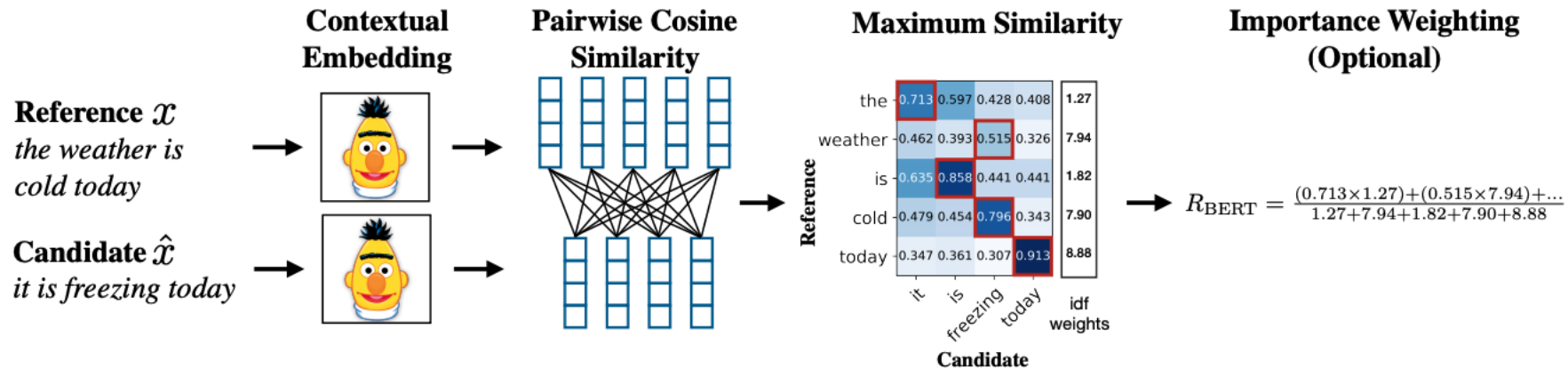


Figure 1: Illustration of the computation of the recall metric R_{BERT} . Given the reference x and candidate \hat{x} , we compute BERT embeddings and pairwise cosine similarity. We highlight the greedy matching in red, and include the optional idf importance weighting.

BERTSCORE The complete score matches each token in x to a token in \hat{x} to compute recall, and each token in \hat{x} to a token in x to compute precision. We use greedy matching to maximize the matching similarity score,² where each token is matched to the most similar token in the other sentence. We combine precision and recall to compute an F1 measure. For a reference x and candidate \hat{x} , the recall, precision, and F1 scores are:

$$R_{\text{BERT}} = \frac{1}{|x|} \sum_{x_i \in x} \max_{\hat{x}_j \in \hat{x}} \mathbf{x}_i^\top \hat{\mathbf{x}}_j, \quad P_{\text{BERT}} = \frac{1}{|\hat{x}|} \sum_{\hat{x}_j \in \hat{x}} \max_{x_i \in x} \mathbf{x}_i^\top \hat{\mathbf{x}}_j, \quad F_{\text{BERT}} = 2 \frac{P_{\text{BERT}} \cdot R_{\text{BERT}}}{P_{\text{BERT}} + R_{\text{BERT}}}.$$

Results : Quality of Generated Queries (1/3)

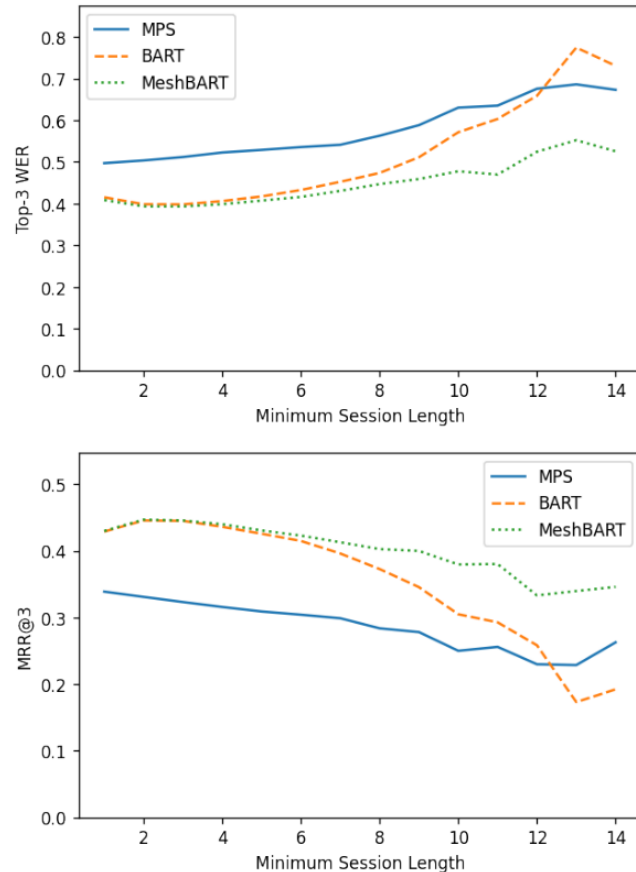


	WER	BertF1	MRR@3	S@3
Seq2Seq+Attn	88.3	53.1	9.8	14.1
MPS	49.7	72.6	33.9	47.1
BART	41.5	76.1	42.5	54.6
MeshBART	40.9	76.5	42.7	55.0

Table 1: Top-3 test performance. Differences between BART and MeshBART on WER and BertF1 are significant ($p < 10^{-4}$) on the Wilcoxon sign-rank test.

- MPS (직전 질의와 Co-occurrence 가장 높은 질의 선택)도 나름 괜찮은 결과
- MPS < BART < MeshBART

Results : Quality of Generated Queries (2/3)



Minimum Session Length
(한 세션에서 입력된 질의 수) 구간 별로 나누어서 보면,

긴 Session 에서 Mesh BART가 잘한다.

-> 긴 Session일수록
Diverse and Complex intents 를
갖고 있기때문에, Next Query 예측이
어렵다.

Figure 2: A breakdown of top-3 word error rate and MRR@3 by minimum session length. Each bucket on x-axis indicates a sub-population of test sessions that contain at least X queries.

Results : Quality of Generated Queries (3/3)



- **MeshBART** vs **MPS** (Generative vs Discriminative)
- W/T/L (Wins, Tie, Loose) Analysis : **MeshBART** vs **MPS**
: (30% Wins / 52% Tie / 18% Loose) on MRR@3
- Sessions with **only one preceding query** (for 39.2% of test session)
 - **MeshBART** can produce at least one **novel** suggestions
(i.e. queries **not seen** in the candidate pool)
(MPS는 후보가 없으면 생성이 불가능)

Preceding Query (Q_n)	Candidates	Generative Suggestions (MeshBART)
environmental technology	environmental, environment	environmental, environmental science, sustainability, environmental scientist
part time adobe	part time, part time marketing	part time marketing, marketing, digital marketing, graphic design
aviation security adelaide	adelaide airport security	aviation security, security, security officer, airport security

*Adelaide
: 호주도시명

Table 2: Query generation examples on **tail queries**, based on test sessions with **only one preceding query** that the logs fail to produce enough candidates for due to scarcity. Generative models such as MeshBART can produce **reasonable suggestions regardless of candidate pool coverage.**

Results : Analysis of Behavioral Hypotheses (1/3)



- Figure 3(a) : 세션의 마지막 질의에 대한 사용자 클릭의 강도
 - > 대부분의 클릭이 **마지막 질의**에서 발생
 - > 사용자가 원하는 결과가 나올 때 까지 질의를 명확히 표현하려고 노력하고, 검색 세션 종료 전에 대부분 검색결과를 소비(클릭)함
- > K4의 가설을 반영
 - K₄ : 사용자의 **직전에 입력한 질의와 반응한 문서**가 중요할 것이다.

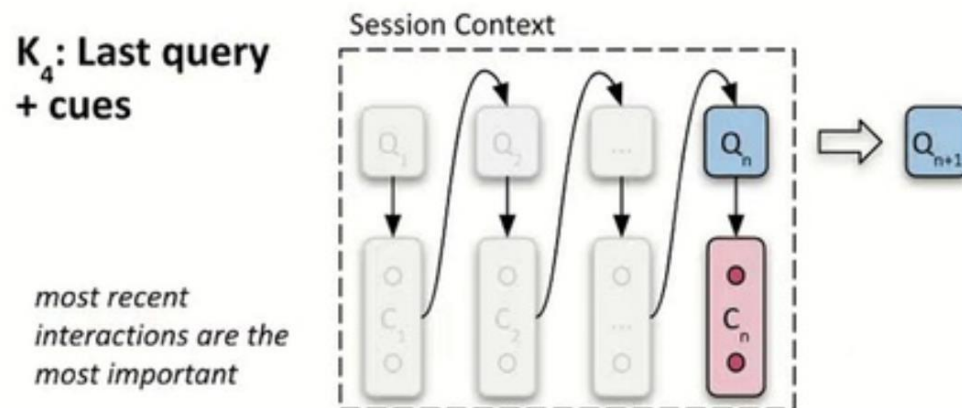
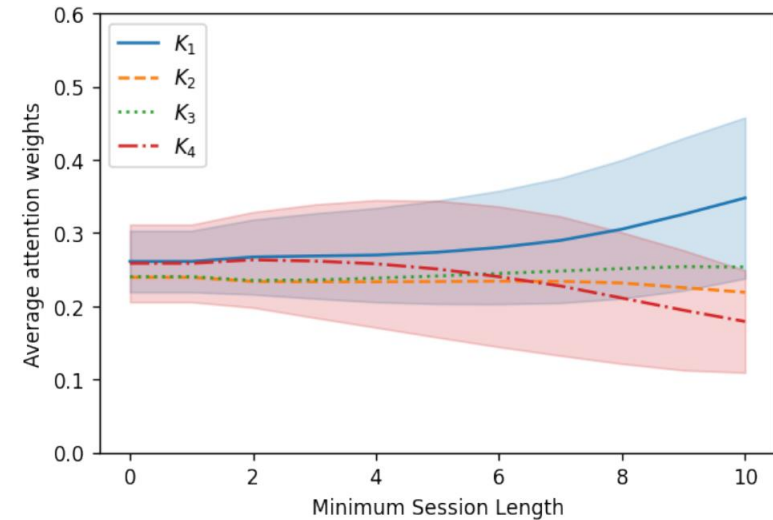
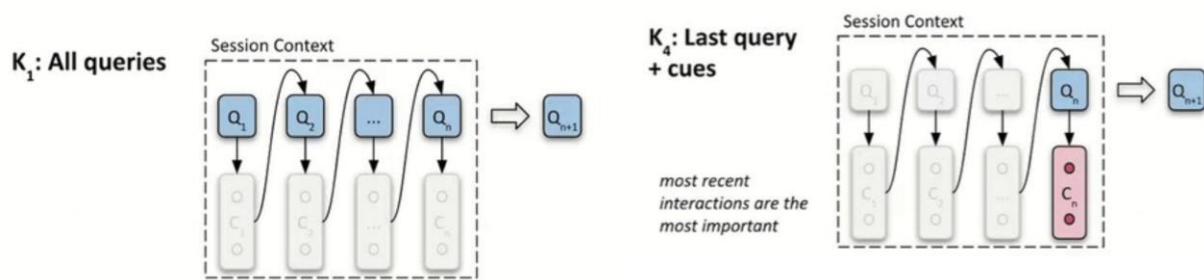


Figure 3: (a) Column-normalized contingency table illustrating clicking behavior.



Results : Analysis of Behavioral Hypotheses (2/3)

- Attention Weight 을 통한 분석
- X: 세션 길이



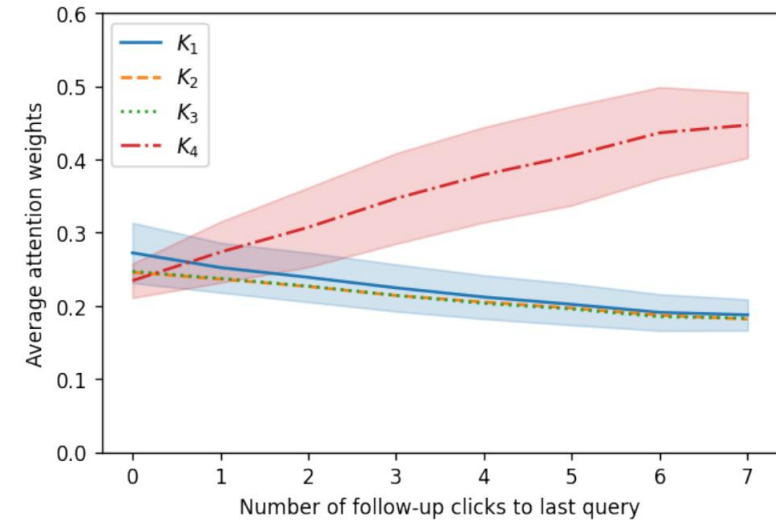
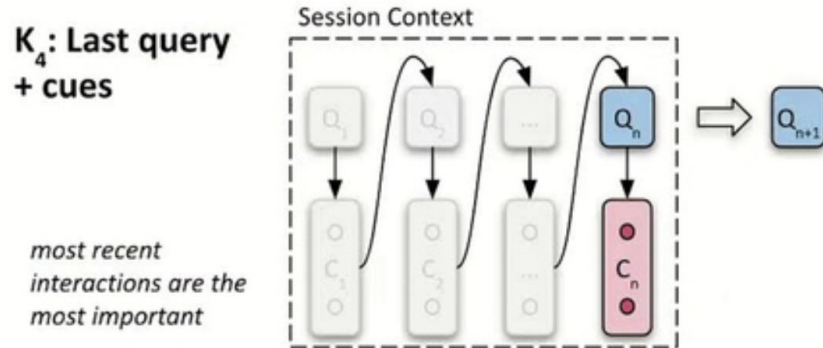
(b) **Average attention weights** for all four behavioral hypotheses
Shaded areas for K1 and K4 indicate the standard error.

- K1(All preceding queries) 의 Attention이 Positively Correlated.
: 더 긴 세션 -> **사용자가 적극적으로 exploring.**
Explicit Search Intent(=Query)가 중요하다는 signal
- K4 는 길어질수록 Attention 이 줄어듦 => Most recent interaction 의 중요도는 떨어짐



Results : Analysis of Behavioral Hypotheses (3/3)

- Attention Weight 을 통한 분석
- X : 마지막 질의에서의 클릭 수



(c) Average attention weights for all four behavioral hypotheses
Shaded areas for K_1 and K_4 indicate the standard error.

- 마지막 질의에서 클릭이 많았던 질의에서 K_4 에 대한 가중치가 높음.
- (b),(c) => flexibility to draw information from different hypotheses in a unified query generation process

Conclusion



- an **effective approach** for **incorporating** user induced interaction patterns as **behavioral hypotheses** into the **query generation** process
- Under an encoder-decoder Transformer framework, the proposed **tokenwise attentions** demonstrate the desirable modeling working by **placing emphasis on different behavioral hypotheses** at different occasions.
- In future work, we will focus on producing novel continuations of the user's search intent, extending the approach to **other domains**, and **automating the design of behavioral hypotheses**.

Reference



- [0] <https://slideslive.com/38939310/incorporating-behavioral-hypotheses-for-query-generation>
- [1] <https://lucidworks.com/ai-powered-search/head-tail-analysis/>
- [2] SORDONI, Alessandro, et al. A hierarchical recurrent encoder-decoder for generative context-aware query suggestion. In: Proceedings of the 24th ACM International on Conference on Information and Knowledge Management. 2015. p. 553-562.
- [3] LEWIS, Mike, et al. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. arXiv preprint arXiv:1910.13461, 2019.
- [4] <https://youtu.be/VmYMnpDLPEo> / 딥러닝논문읽기모임 : BART paper review
- [5] <https://jalammr.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>
- [6] <https://hcnoh.github.io/2018-12-11-bahdanau-attention>
- [7] BAHDANAU, Dzmitry; CHO, Kyunghyun; BENGIO, Yoshua. Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473, 2014.
- [8] https://en.wikipedia.org/wiki/Mean_reciprocal_rank
- [9] ZHANG, Tianyi, et al. Bertscore: Evaluating text generation with bert. arXiv preprint arXiv:1904.09675, 2019.