

1 We thanks reviewers for their feedback. First, we address some common concerns.

2 **Disparity between the region caption based queries in our experiments and the user queries in real scenario.**

3 As emphasized in our manuscript (L53), leveraging region captions for weak supervision can be seen as one of the
4 advantages of our method, as no annotation is required. While we agree strong supervisory signals such as real user
5 queries could bridge the domain gap and would like to explore further in this direction, we choose at this stage to use
6 only "weak but free" signals to see to what extent they can be generalized to practical applications. We demonstrate
7 the effectiveness of the proposed method on simulated data with extensive experiments, on real scenarios via human
8 subject studies on Amazon Mechanical Turk. Our experiments show the region caption based queries can potentially be
9 generalized to real scenarios with high promise.

10 **Patterns of real user queries**

11 Compared with the simulated caption data, the real user queries we collected show distinct patterns, e.g. long
12 descriptions of the target scenes, less informative sentences (e.g. 'good!'). That's why we observe smaller performance
13 gaps between our method and the alternative approaches in the user study. We will include more details of the user
14 study and examples of the real user queries in the revised paper.

15 **Reviewer 1**

16 **#Q2: Baselines**

17 With relatively less prior works on this research, we compare the proposed method with variants of state-of-the-art
18 approaches for the most related topics, e.g. dialog-based interactive product search. We will incorporate the suggested
19 baselines (e.g. late score/rank fusion vs early feature fusion, linear query encoding vs hierarchical query encoding) in
20 the revised paper.

21 **#Q4: Applications of the proposed method**

22 We envision the proposed method could generally help with natural image search. Potential applications include
23 retrieving very specific images of complex scenes the users encountered before, or exploring inspiring images for
24 creative content generation (e.g. Adobe Stock Image).

25 **Reviewer 2**

26 **#Q1: Incorporating advanced language models**

27 We have explored using bidirectional language encoders and found it performs similar with unidirectional encoders in
28 this task. We conjecture that unidirectional and bidirectional encoders provide comparable contextual signals when
29 encoding the per-turn query as a single feature vector for downstream modules. In the current manuscript, we focus
30 more on the sequential encoding of multiple sentences, and would like to explore and incorporate more advanced
31 language models such as BERT in the future.

32 **Reviewer 3**

33 **#Q1: Distinguishing the paper's contribution with memory networks**

34 In contrast to the previous sentence encoding methods which perform **query** and possibly **update** operations on a
35 predefined external memory space (e.g. the agenda items in Kiddon et al. 2016, neural checklist models), we focus on a
36 more challenging scenario where the model needs to **create** and **update** the memory module (the state vectors in our
37 case) **on-the-fly** so as to maintain the dynamic states of multiple-turn queries. We will elaborate more to distinguish our
38 method with memory networks in the revised paper.

39 **#Q2: Experimental details**

40 (1) The region captions and their orders are randomly sampled. We keep the captions and their orders of the validation
41 and test sets unchanged for all our experiments; (2) We use ten turns in all our simulated experiments as we'd like to
42 track and demonstrate the performance of the proposed method in both short-term and long-term scenarios, as shown in
43 Fig. 3. In the user study, we start with ten-turn queries but observe the users are less willing to continue and finish the
44 tasks if they could not succeed in five turns, so we evaluate the five-turn queries in our experiment; (3) We use different
45 image sets for training, validation and evaluation (L208), where the images retrieved are from the corresponding sets at
46 different stages respectively. All the evaluations (including the user study) are performed on the test set, which contains
47 9896 images (L208); (4) All images in the candidate set are ranked in each turn. (5) Faster RCNN is NOT finetuned in
48 our experiments.

49 **#Q3: Number of the placeholders**

50 We confirm in our experiments that more placeholders result in better performance but also lead to scalability difficulties.
51 One of the main goal of the proposed work is to model multiple-turn queries with dynamic lengths using a fixed set
52 of hidden states (fixed computational budget accordingly). We're also happy to include the suggested experiments and
53 improve the presentation of the paper in the revised version.

54 **#Q4: Questions about the human evaluation**

55 As answered in #Q2, we start experimenting with ten queries but discover it cannot fit in the real scenario. We agree
56 that using less placeholders will be more convincing in this case and will rerun the experiment in the revised paper.