

1 We thank the reviewers for their thoughtful comments.

## 2 Reviewer 1

3 a) We have conducted preliminary model dynamics and error analysis based on your feedback and summarize our  
4 results below. We will add them to an updated version of the paper and perform more analysis to improve it. This  
5 analysis was done on the question answering experiment.

6 **Examples where local adaptation helps.** In Table 1, we show two test examples where our model answers  
7 incorrectly before local adaptation, but correctly after. In the first case, we can see that training examples retrieved  
8 from memory are thematically related to the test example. In the second case, since the query is shorter, retrieved  
9 training examples tend to be more syntactically related. Although we only show the two nearest neighbors for each  
10 query here, our analysis provides an insight on ways our model uses its memory to improve predictions.

Table 1: Two examples where local adaptation helps.

---

**Context:** david niven ( actor ) - pics , videos , dating , & news david niven male born mar 1 , 1910 james  
david graham niven , known professionally as david niven , was an english actor and novelist [ . . . ]

---

**Query:** in 1959 , for which film did david niven win his only academy award ?

---

**First two training examples retrieved from memory (2 nearest neighbors):**

in which of her films did shirley temple sing animal crackers in my soup ?

in 1968 , which american artist was shot and wounded by valerie solanis , an actress in one of his films ?

---

**Context:** dj kool herc developed the style that was the blueprint for hip hop music . herc used the record to  
focus on a short , heavily percussive part in it : the " break " . [ . . . ]

---

**Query:** what was the break ?

---

**First two training examples retrieved from memory (2 nearest neighbors):**

what was the result ?

what was the aftermath ?

---

11 **Relevant examples that are difficult to retrieve.** In Table 2, we show two relevant training examples (as judged  
12 by humans) that are difficult to retrieve by the model (they are not in the 1,000 nearest neighbors) for the query  
13 what was the name of bohemond s nephew. The two relevant training examples ask about the nephew of a  
14 person, which is relevant for the given query. However, since they are phrased differently to the query, they are far in  
15 the embedding space, which is why a nearest neighbor method fails to retrieve these training examples. Our analysis  
shows that a better embedding and/or retrieval method can potentially improve the performance of our model.

Table 2: Relevant examples that are difficult to retrieve from memory.

---

**Query:** what was the name of bohemond s nephew

---

**Relevant examples not retrieved (Euclidean distances to the query in parentheses):**

(87.88) who was the nephew of leopold

(103.96) who is the nephew of buda king casimer iii the great

---

16  
17 b) We will reorganize the presentation of the main results (Table 1 in the paper) to include some results from Appendix  
18 B such that they are more informative for readers who work with datasets we consider in our paper.

19 c) We will revise the characterization of our work with respect to McCann et al. according to your suggestion.

## 20 Reviewer 2

21 a) Examples from each dataset need not be seen contiguously for our model to work. In the limit, when all examples  
22 across datasets are shuffled, we reach the performance of the multitask upper bound shown in the paper.

23 b) For classification, yes, it is correct that we have a 33-way softmax. These 33 classes already include overlapping  
24 classes from two datasets (Yelp and Amazon).

## 25 Reviewer 3

26 a) Thank you for your suggestions and pointers to typos. We will add more details about our experiments with other  
27 reading and writing strategies as an appendix.