

1 **Reviewer 1**

2 > *“The results (Table 2) show similar accuracies for all attention models.”*

3 Note that the VQA results in Table 2 with continuous attention use fewer basis functions than discrete regions. Although
4 the accuracies are similar, the unimodal attention suggests better interpretability (as noted by R2 and R3).

5 > *“It would help to have a short algorithm describing how to implement the forward and backward passes efficiently.”*

6 Good idea, we will add this to the camera-ready version.

7 > *“Line 134: “a condition (...) is g being strongly convex”. Is this a necessary or a sufficient condition?”*

8 Sufficient; we will clarify and follow the suggestions (move the beta-escort definition to the main text and fix typos).

9 **Reviewer 2**

10 Thanks for the positive comments and for pointing out the work of Cordonnier et al. (2020). We will add a citation.

11 > *“Have you experimented with other ways such as linear interpolation (1D) or bilinear (2D)?”*

12 We chose ridge regression as it enables a closed-form solution expressed linearly in terms of the basis functions (Eq. 15)
13 and matrix G can be precomputed, leading to a fast implementation. Note that the proximity of tokens/pixels is taken
14 into account (the basis vectors $\psi(t_\ell)$ forming F are located at each token/pixel). We haven’t tried linear interpolation,
15 but this is an interesting suggestion (although it might make attention computation more challenging).

16 > *“Does it remove the need for additive positional encoding?”*

17 Very good point; this is indeed one advantage of our approach – by converting the input to a function on a predefined
18 continuous space, it encodes “positions” implicitly in a natural way, not requiring explicit positional encoding.

19 **Reviewer 3**

20 > *“The proposed method’s application on VQA is limited to grid feature.”*

21 Actually, our method can handle BUTD features too: it suffices to let the t_ℓ coordinates in the multivariate regression
22 (Eq. 15) be placed on these regions instead of on a grid. However, we opted not to rely on an external object detector,
23 in order to check if continuous attention has the ability to detect relevant objects on its own (see ellipses in Fig. 3).
24 However, for a high-level vision system, combining our method with BUTD is an interesting idea.

25 > *“Why model text as continuous inputs? Text are naturally discrete tokens.”*

26 We agree text is fundamentally a discrete sequence of symbols. However, when processing long documents or attending
27 to snippets, modelling it as a continuous signal may be advantageous, due to smoothness and independence of length.

28 > *“The proposed method assumes attention probability is single mode. Is this a reasonable assumption?”*

29 Good point. Unimodal attention is useful to focus on a single object or text segment of varying size, avoiding
30 “fragmenting” attention probability; however, in some applications, multimodal attention may be preferable. Our method
31 can be extended to multiple modes via a suitable choice of $\phi(t)$ (e.g., a mixture of Gaussians), but this will require
32 numeric integration for attention computation. A simpler strategy (see lines 258-260) is to use multi-head or sequential
33 attention.

34 > *“Can it be applied to deep transformer models with multi-head attention?”*

35 Great question. We have ongoing work applying this to transformer models (but out of scope for this paper). Briefly,
36 the computation cost is $O(N)$ for each attention head (against $O(L)$ in the discrete case) where $N \ll L$ is the number
37 of basis functions, plus an extra $O(NL)$ cost in the first layer to perform the multivariate regression on L tokens.

38 > *“Is there any intuition (...) how the continuous attention can improve the accuracy of downstream tasks?”*

39 In general, continuous attention can make it easier to attend to large spaces with different resolution levels, with a fixed
40 number of Gaussian RBFs with several variances. It can also lead to more focused attention (the VQA experiments
41 suggest this) and better control of time steps with continuous data streams (e.g., irregularly sampled time series). We
42 haven’t explored all these directions, but we believe these are promising areas of future research.

43 **Reviewer 4**

44 Please check our answers to R1 and R3 above (the answer is yes to positional encodings).