

# Probabilistic Multimodal Learning with Bayesian Disagreements

## Supplementary Material

In the supplementary section, we provide additional details about the method. We also provide additional details about the results.

### 1. Additional Details about the Methodology

#### 1.1. Attention Reuse and Partial Sum algorithm

The algorithm to describe the Attention Reuse and Partial Sum Mechanism (Section 3.1.2) is described in Algorithm 1.

---

#### Algorithm 1 Attention Reuse and Partial Sum

---

**Require:**  $Q, K, V \in \mathbb{R}^{n \times b \times d}$ , number of modalities  $M_n$

1: **Given:**  $Q \in \mathbb{R}^{n \times b \times d}$ ,  $K \in \mathbb{R}^{n \times b \times d}$ ,  $V \in \mathbb{R}^{n \times b \times d}$   
 $\{n$ : batch-size,  $b$ : number of tokens,  $d$ : feature dimension};  $M_n$  is the number of modalities

2: **Define:**

3:  $M_I = \{0, 1, \dots, M_n - 1\}$

4:  $I_M = \{I_1, I_2, \dots, I_{M_n}\}$   $\{I_i$  represents the indices of tokens for  $i$ -th modality}

5:  $\mathcal{P} \leftarrow \mathcal{P}(I_M) \setminus \{\emptyset\}$  {Power set of  $I_M$ , excluding the empty set}

6:  $\mathcal{P}_M \leftarrow \mathcal{P}(M_I) \setminus \{\emptyset\}$  {Power set of  $M_I$ , excluding the empty set}

7: **Step 1: Full Attention Calculation**

8:  $\mathcal{A} \leftarrow \exp\left(\frac{Q \cdot K^T}{\sqrt{d}}\right)$

9: **Step 2: Initialize containers**

10:  $M \leftarrow []$ ;  $E \leftarrow []$ ;  $O \leftarrow []$  {Array to store output of attention}

11: **Step 3: Compute intermediate attention outputs**

12: **for** each  $i \in \{0, 1, \dots, M_n - 1\}$  **do**

13:  $M[i] \leftarrow \mathcal{A}[:, I_M[i]] \cdot V[I_M[i], :]$

14:  $E[i] \leftarrow \sum_{j \in I_M[i]} \mathcal{A}[:, j]$

15: **end for**

16: **Step 4: Reuse attention and compute output**

17: **for** each  $(I, J) \in (\mathcal{P}, \mathcal{P}_M)$  **do**

18:  $O.append\left(\frac{\sum_{i \in J} M[i]}{\sum_{i \in J} E[i]}\right)$  {Row-wise division}

19:  $O[-1] \leftarrow O[-1][I, :]$  {Slice}

20: **end for**

21: **Note:**  $\cdot$  operator represents matrix multiplication.

---

#### 1.2. Probabilistic Attention

A visualization to show the difference between the Traditional Attention block and the proposed Probabilistic Attention block is shown in Figure 1. Figure 1(a) shows

the traditional attention block and Figure 1(b) shows the proposed probabilistic attention block.

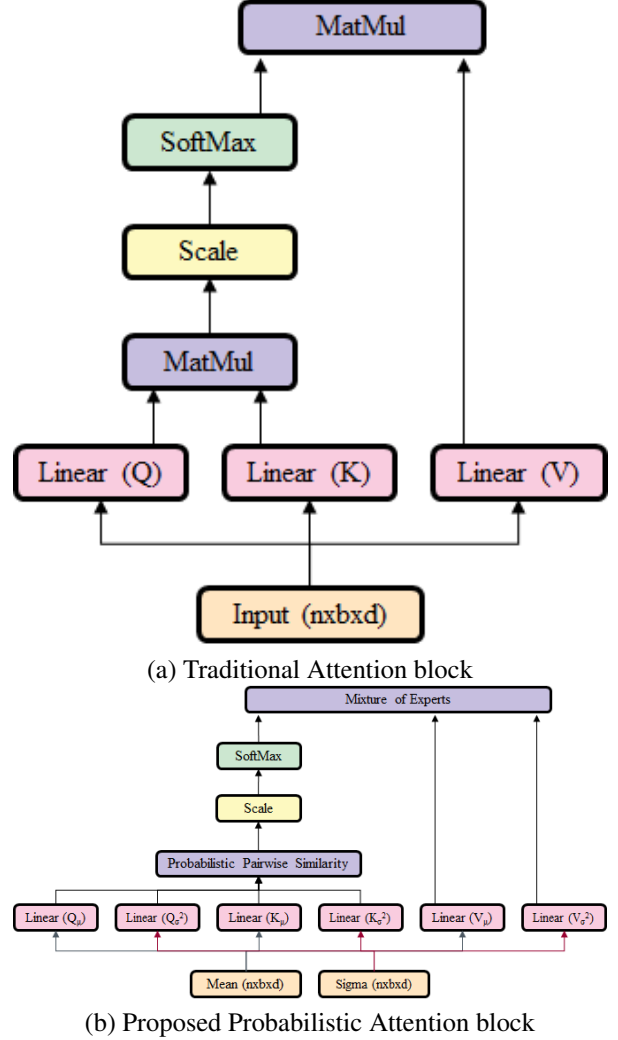
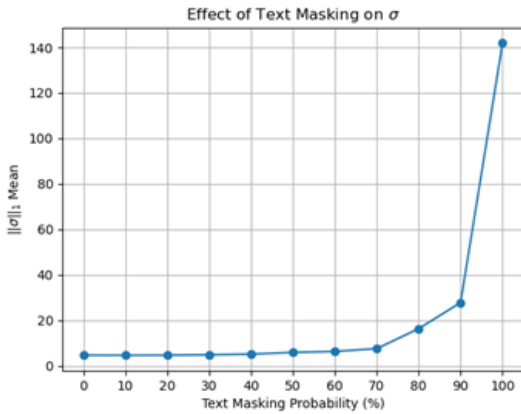


Figure 1. Comparison of our proposed probabilistic attention block and traditional attention.  $n$  is batch size,  $b$  is the number of tokens and  $d$  is the dimensionality.

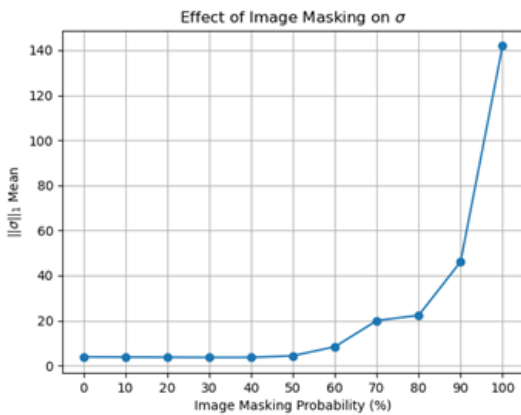
### 2. Additional Details about the Results

#### 2.1. Effects of $\sigma$

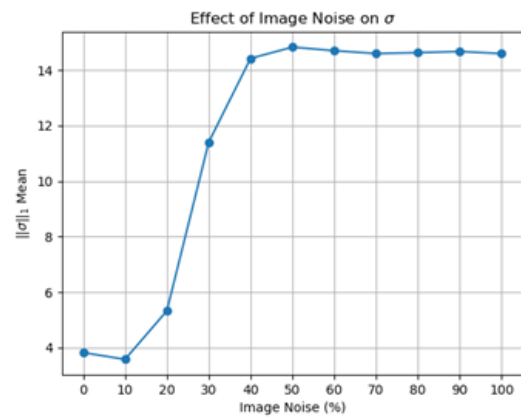
We analyze the sigma (standard deviation) value with increasing amounts of degradation in the input data (Figure 2). We experiment with three settings. First is increasing the text masking ratio (Figure 2(a)). Second is increasing the image masking ratio (Figure 2(b)). Lastly, we increasingly add Gaussian noise (Figure 2(c)). The



(a) Increasing text masking ratio



(b) Increasing image masking ratio



(c) Increasing noise in image

Figure 2. Sigma values captured by our PMBD model under various amounts of degradation severity.

X-axis has increasing amounts of degradation. In all cases, we can see that our model outputs sigma values that have approximately a monotonic behavior with the degradation severity. This implies that as the reliability

decreases (that is, the degradation increases), the predicted sigma value also increases.

We also qualitatively visualize the sigma values outputted by our model using the PMBD-PA model (Figure 3). We can see that the informative regions of the image have lower sigma values. Lower sigma values are indicated by blue color. The text regions in the images usually show lower sigma values indicating a high information region.



Figure 3. Qualitative visualizations of sigma values for various images in the MM-IMDB dataset.

## 2.2. Results on Epic-Kitchens dataset

In the paper in Figure 5, we show the results of different methods on the Action recognition task under noisy conditions. However, we only show the results on the

Action score. In this Supplementary Material, in Table 1, we show the results on the Verb, Noun and Action recognition under noisy conditions to validate that the trends in Action are similar to the trends in Verb and Noun recognition.

Table 1. Accuracy on Epic-Kitchens dataset. ‡ refers to supervised fine-tuning.

<b>Method</b>	<b>Verb</b>	<b>Noun</b>	<b>Action</b>
<i>Background Noise SNR – 30</i>			
CoMM‡	0.651	0.515	0.379
ProbCoMM‡	0.688	0.523	0.380
<b>PMBD‡</b>	0.711	0.568	0.403
<b>PMBD-PA‡</b>	0.723	0.577	0.412
<i>Background Noise SNR – 20</i>			
CoMM‡	0.512	0.451	0.281
ProbCoMM‡	0.515	0.459	0.293
<b>PMBD‡</b>	0.632	0.504	0.341
<b>PMBD-PA‡</b>	0.641	0.509	0.350
<i>Background Noise SNR – 15</i>			
CoMM‡	0.501	0.443	0.272
ProbCoMM‡	0.504	0.449	0.284
<b>PMBD‡</b>	0.628	0.500	0.336
<b>PMBD-PA‡</b>	0.638	0.502	0.347
<i>Background Noise SNR – 10</i>			
CoMM‡	0.496	0.439	0.267
ProbCoMM‡	0.498	0.440	0.278
<b>PMBD‡</b>	0.625	0.498	0.333
<b>PMBD-PA‡</b>	0.636	0.499	0.345
<i>Gaussian Blur in Video</i>			
CoMM‡	0.421	0.415	0.243
ProbCoMM‡	0.452	0.429	0.256
<b>PMBD‡</b>	0.495	0.478	0.272
<b>PMBD-PA‡</b>	0.500	0.486	0.280