# Supplementary Material for Unpaired Image-to-Speech Synthesis with Multimodal Information Bottleneck

Shuang Ma	Daniel McDuff	Yale Song
SUNY Buffalo	Microsoft	Microsoft
Buffalo, NY	Redmond, WA	Redmond, WA
shuangma@buffalo.edu	damcduff@microsoft.com	yalesong@microsoft.com

### 1. Network Architectures and Parameter Settings

We provide implementation details of our model with the parameter settings used in our experiments. We encourage the readers to refer to Figure 3 and Figure 4 of our main paper when reading this section. We use the following notations to refer to commonly used computation blocks in the neural networks: Conv1D(#channels, kernel\_size, stride\_size), Conv2D(#channels, kernel\_size, stride\_size), FC(#units), GRU(#units).  $\oplus(\cdot)_{res}$  refers to the residual connection. We use the superscript *j* to refer to the modalities  $j \in \{img, txt, spch\}$ .

### 1.1. Encoders (Figure 3 (left) in the main paper)

- Image encoder:  $\mathbf{x}^{img} \rightarrow \text{Conv2D}(64, 4, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv2D}(128, 4, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv2D}(256, 4, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv2D}(512, 4, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{MaxPool} \rightarrow \mathbf{e}^{img}$
- Text encoder:  $\mathbf{x}^{txt} \rightarrow \text{LookupTable(66, 128)} \rightarrow \text{FC}(256) \rightarrow \text{ReLU} \rightarrow \text{Dropout}(0.5) \rightarrow \text{FC}(128) \rightarrow \text{ReLU} \rightarrow \text{Dropout}(0.5) \rightarrow \text{CBHG } [2] \rightarrow \text{AvgPool} \rightarrow \text{FC}(512) \rightarrow \tanh \rightarrow \mathbf{e}^{txt} \in \mathbb{R}^{512}$
- Speech encoder:  $\mathbf{x}^{spch} \rightarrow \text{Conv2D}(32, 3, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv2D}(64, 3, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{GRU}(256) \rightarrow \text{FC}(512) \rightarrow \text{tanh} \rightarrow \mathbf{e}^{spch}$

### 1.2. Multimodal Information Bottleneck (Figure 4 in the main paper)

• Modality transformer:  $e^j \rightarrow FC(256) \rightarrow ReLU \rightarrow \oplus (Conv1D(256, 1, 1) \rightarrow ReLU \rightarrow BN \rightarrow Conv1D(256, 1, 1) \rightarrow Conv1D(256, 1) \rightarrow Co$ 

 $\operatorname{ReLU} \to \operatorname{BN} \left( \right)_{res} \to \operatorname{tanh} \to \mathbf{z}^{j}$ 

- Memory fusion module:
  - 1. Define: Memory  $M \in \mathbb{R}^{n_k \times d_k/n_{heads}}$ , where  $n_k = 128, d_k = 256, n_{heads} = 4$
  - 2. Query  $\mathbf{q}^j \mathbf{z}^j$ , Key k Conv1D(256, 1, 1)(tanh(M)), Value v tanh(M)
  - 3.  $(\mathbf{q}_h^j, \mathbf{k}_h, \mathbf{v}_h)$ SplitHeads $(\mathbf{q}^j, \mathbf{k}, \mathbf{v}), h = 1, \cdots, n_{heads}$
  - 4.  $\alpha_h^j \text{SoftMax}\left(\mathbf{q}_h^j \mathbf{k}_h / \sqrt{d_k}\right), h = 1, \cdots, n_{heads}$
  - 5.  $\mathbf{u}_h^j \alpha_h^j \times \mathbf{v}_h, h = 1, \cdots, n_{heads}$
  - 6.  $\mathbf{u}^j$  ConcatHeads $(\mathbf{u}_h^j)$

### 1.3. Decoders (Figure 3 (right) in the main paper)

- Image decoder:  $\mathbf{u}^{txt} \rightarrow \text{Conv2D}^{\intercal}(32, 4, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv2D}^{\intercal}(16, 4, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv2D}^{\intercal}(8, 4, 2) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv2D}^{\intercal}(8, 4, 2) \rightarrow \text{tanh} \rightarrow \mathbf{y}^{img}$
- Text decoder:  $\mathbf{u}^{img} \rightarrow \text{Dropout}(\text{LSTM}(128), 0.3) \rightarrow \text{Dropout}(\text{LSTM}(128), 0.3) \rightarrow \text{FC}(n_{symbols}) \rightarrow \text{SoftMax} \rightarrow \mathbf{y}^{txt}$
- Speech decoder: u<sup>txt</sup> → AttentionRNN(GRU 256) → DecoderRNN6(Dropout(LSTM(256), 0.3) → Dropout(LSTM(256), 0.3)) → reshape → y<sup>mel-spectrogram</sup> → CBHG [2] (80 mels) → FC(1025) → Griffin-Lim (y<sup>linear</sup>) → y<sup>speech</sup>



Piecewise: A room with blue walls Ours: A room with walls that

are painted blue.



Piecewise: A man is walking with some sheep. Ours: A person is walking with some sheep.



Piecewise: A big bed in a shady room Ours: A big bed in a shady bedroom with two black bags.



Piecewise: A dog looking outside Ours: A dog looking outside.



Piecewise: A table with various colored flower vases. Ours: A white table with five colored flower vases.



Piecewise: A train engine on a track Ours: A red train engine on the

track under a bridge.

are on a white plate with chili.

**Ours**: Some breads with bowls

**Piecewise**: A boy with a hamet

standing with another boy.

Ours: Two boys with a

skateboard and helmet

standing in the dark.

of red sauce nearby.



Piecewise: A teddy bear lays in bea **Ours**: A teddy bear lays in bed.



in a chair. Ours: A cat is sitting in a chair.



Piecewise: Some grilled cheese Piecewise: A truck and a car parked in a driveway. Ours: A green car and a white truck parked in a driveway.



flowers are in a blue vase.

a blue vase.



racquet

Piecewise: A woman is playing tennis Ours: A woman is standing in some trees with a tennis



Piecewise: A giraffe grazing in a field Ours: Two giraffes standing outside in the field.



Piecewise: A bowl with chips in the kitchen. Ours: A bowl full of food sitting on a kitchen counter.



Piecewise: Some people are all sittina toaether. Ours: Some kids are sitting together with a teddy bear.



Piecewise: A building with a tower and red wall. Ours: A building with a tower and a clock on the wall.

Figure 1. Image-to-speech synthesis results. Green: Fine-grained and correct instances synthesized by our model. Red: incorrect pronunciation synthesized by the piecewise model. Audio samples are available at https://bit.ly/2U7741S

Ours: Pink and white flowers in

### 2. Skip-Modal Synthesis Results

Figure 1 shows additional image-to-speech synthesis results; we manually transcribed the synthesized audio outputs for the purpose of presentation. Consistent with the qualitative results reported in the main paper (Figure 5), we see that our approach produces more detailed descriptions and has a larger vocabulary than the baseline. We encourage the readers to visit our anonymized website and listen to the audio samples: https://bit.ly/2U7741S

### 3. Cross-Modal Retrieval Results

Besides the synthesis tasks, another way to evaluate the performance of our model is via cross-modal retrieval. In this section, we show qualitative results of cross-modal retrieval where we use an instance from either dataset and find the most

Piecewise: A cat that is sitting



Piecewise: Three men

standing in the grass.

outside in the grass.

Ours: Three men standing

Piecewise: Two men are

Ours: Two people are talking

talking in the room.

in the kitchen.

Piecewise: A man with a surfhoard Ours: A man carrying a surfboard.

#### Image Query



Speech Query (EMT-4):

Yeah, a giant zebra.

I'm jealous of a dog.

Text (COCO) Top1: A cat sitting on top of a window ledge in the daytime. Top2: A fat orange cat sitting on a deck. Top3: Two cats sitting on top of a window sill.

#### Speech (EMT-4)

Speech (EMT-4)

Top1: I choose cats all the way. Top2: I miss all my cats! Top3: I felt so awful; leaving my cats behind.

Top2: And I want pasta or sashimi.

Text (COCO)
Top1: A plate of pasta with meat and broccoli
together.
Top2: There is a large plate of pasta and
broccoli.
Top3: A plate of pasta with broccoli, beef, and
cheese.

Top1: A zebra has its neck turned and is looking

Top3: Two zebras that are standing in the grass.

Top2: A zebra standing with its reflection in a

Top1: There is a log playing with a dog bed. Top2: A dog jumping in the air with a Frisbee in

Top3: A black and white dog stands beside a

Text (COCO)

on its side.

Text (COCO)

its mouth.

person.

pool.

# Image (COCO)

bass..



Top1: I was eating pasta and I bit the fork and it broke.

Top3: The irrepressible chef, Mollie Ahlstrand, is a genius with homemade breads, fresh pasta and see



<b>Text Query (COCO):</b> Group of baseball players playing on a wet field.	<b>Speech (EMT-4)</b> Top1: <i>I miss baseball.</i> Top2: Sorry! <i>I can't stand baseball</i> ! <i>I can't stand you.</i> Top3: We start with baseball with the L.A Dodgers.	Image (COCO)
A birthday cake shaped	Speech (EMT-4)	

like a wrestling ring.

Top1: Hope it is a great birthday. Top2: Happy birthday by the way. Top3: Celebrating a birthday on night shift.





Figure 2. Cross-Modal retrieval results. The first column shows queries from each modality. The second and third columns show the top-3 retrieval results from the other two modalities. Audio samples are available at https://bit.ly/2U7741S

similar instances from different modalities from both datasets. Specifically, we compute  $\mathbf{u}^{j}$  from all instances in the test splits of both datasets, and compute the cosine similarity between any pair of cross-modal instances.

Figure 2 shows the top 3 retrieved results in different combinations of modalities. We can see that the retrieved results are very related to the query at the object level, e.g., "dog" and "zebra" in the first and the second rows, while on the other two rows the results are related to the query at the scene/context level, e.g., "baseball game" and "birthday party". It is particularly interesting to see that the results are reasonable even for the cross-dataset retrieval settings (using an image from COCO to retrieve audio/speech instances from EMT-4). This suggests the representations extracted using our model are not very sensitive to the dataset and the modalities involved.

Batch Sam-	B@1	WER	WER
pling Strategy	(I2T)	(S2T)	(T2S)
Alternative	74.1	3.88	10.5
Mixing	74.5	3.76	10.5

Table 1. Evaluation of different batch sampling strategies. I2T: image-to-text, S2T: speech-to-text, T2S: text-to-speech.

$n_k (\text{fix } d_k = 256)$	10	128	256
BLEU-1	62.5	74.1	73.9
$d_k \text{ (fix } n_k = 128)$	64	128	256
BLEU-1	49.3	70.2	74.1

Table 2. Sensitivity of the memory fusion module.  $n_k$ : the number of basis vectors,  $d_k$ : the size of each basis vector.

#### 4. Additional Ablation Experiments

### 4.1. Different Batch Sampling Strategies

As we trained our model on a combination of two datasets, there comes two ways to perform mini-batch training: one that samples instances from only one dataset and alternates between the two (alternate); and another that always samples instances from both datasets (mixing). We compare these two batch sampling strategies in this section. Specifically, in the first setting (alternative) we sample eight instances from either the COCO or EMT-4 dataset, while in the second setting (mixing) we sample four instances from COCO and the other four from EMT-4. We evaluate this on image-to-text (I2T), speech-to-text (S2T), and text-to-speech (T2S) synthesis tasks, reporting BLEU-1 for I2T and the word error rate (WER) for the other two.

Table 1 shows that the performance improves when we use the mixed batch sampling strategy. The improvement is especially pronounced for the text-sensitive tasks; on image-to-text synthesis the BLEU-1 is improved from 74.1 to 74.5, and on speech-to-text synthesis the WER is reduced from 3.88 to 3.76. We did not find significant differences in the text-to-speech synthesis task.

## 4.2. Sensitivity Analysis of Memory Fusion Module with parameters $n_k$ and $d_k$

As we showed in Table 4 (ablation results) in the main paper, the memory fusion module plays an important role in our model; the performance drops most significantly when we bypass this module. It extracts compact, modality-agnostic representations of the multimodal inputs following the information bottleneck principle [1], using the shared external memory M to "bottleneck" any redundant and modality-specific information from leaking into the output representation. To better understand the behavior of this module, we analyze the sensitivity of the module to two hyper-parameters: the number of basis vectors  $(n_k)$  and the size of each basis vector  $(d_k)$  inside the external memory variable M. We evaluate this on image-to-text generation (i.e., image captioning) and report the results using BLEU-1 as our metric.

Table 2 shows our model is more sensitive to the dimension of each basic vector  $d_k$  than the number of basis vectors  $n_k$ ; it achieves a significantly lower performance with  $d_k = 64$  compared to any other combination of the two parameter values. The performance improves as we increase  $d_k$ , achieving the best performance when  $d_k = 256$ ; we did not evaluate beyond  $d_k = 256$  due to the limitations on the GPU memory. As for the number of basis vectors  $n_k$ , we can see the performance is low when there are only a few of them ( $n_k = 10$ ). This shows we need a large number of basis vectors to capture the variety of information contained in multimodal data. We found that the performance is relatively stable when  $n_k$  is greater than 128.

### References

- [1] Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. In *Allerton Conference on Communication, Control, and Computing*, 1999.
- [2] Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly, Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. Tacotron: Towards end-to-end speech synthesis. In *INTERSPEECH*, 2017.