FLIPDIAL: A Generative Model for Two-Way Visual Dialogue Supplementary Material

Daniela Massiceti University of Oxford, UK daniela@robots.ox.ac.uk N. Siddharth University of Oxford, UK nsid@robots.ox.ac.uk Puneet K. Dokania University of Oxford, UK puneet@robots.ox.ac.uk

Philip H.S. Torr University of Oxford, UK

phst@robots.ox.ac.uk

A. Glossary

- **block dialogue/architecture** Models $\mathbf{B}/\mathbf{B}_{\mathbf{AR}}$ are built and trained for the task of two-way visual dialogue (2VD) with data x = d and condition variable $y = \{i, c\}$. Since d refers to the whole dialogue sequence/block $\langle (q_t, a_t) \rangle_{t=1}^T$ we refer to $\mathbf{B}/\mathbf{B}_{\mathbf{AR}}$ as *block* architectures.
- **generation** This represents the scenario when only the condition variable y is available at test time. In this case, the decoder network receives a sample $z \sim p_{\theta}(z \mid y)$, a multivariate Gaussian parametrised by μ_p and exponentiated $\log \sigma_p^2$ learned using the *prior* network. We call the decoded output \hat{d} a generation.
- **reconstruction** Differing from a generation, both y and d are available. The decoder network receives a sample $z \sim q_{\phi}(z \mid d, y)$, a multivariate Gaussian parametrised by μ_q and exponentiated log σ_q^2 learned using the *encoder* network. We call the decoded output \hat{d} a *reconstruction*. The reconstruction pipeline is used during training when the input d and the condition variable y are available. Note, this pipeline is also used when $\mathbf{B}/\mathbf{B}_{AR}$ are evaluated iteratively (see §4.2).

B. Extended Quantitative Results on 1vD task

Tab. 3 in the main paper evaluates **A** and **B**/**B**_{AR} in the task of one-way visual dialogue (1VD). Here we shed light on these numbers and the metrics used to obtain them. We also present a more extensive quantitative analysis of **B**/**B**_{AR} in the 1VD task (see Tab. 6).

Evaluating B/B_{AR} on 1VD We extend Tab. 3 with-Tab. 6, which further compares B/B_{AR} under the iterative evaluation settings of d-qa and $d-q\hat{a}$, using the cross entropy (CE) and Kullback-Leibler (KL) terms of the evidence lower bound (ELBO) and our two new metrics, $\sin_{c,q}$ and $\sin_{(1)}$. We observe that $\mathbf{B}/\mathbf{B}_{AR}$ (*d*-*qa*) shows superior performance of around 7-10 points in MR over $\mathbf{B}/\mathbf{B}_{AR}$ $(d-q\hat{a})$, and also improves in MRR and recall rates. This is expected since the ground-truth rather than predicted answers are included in the dialogue history (along with the ground-truth questions). The metrics $sim_{cap,q}$ and sim_{\circlearrowleft} , on the other hand, show very little performance difference across the two evaluation settings. We also note that ranking performance is worse when both image i and caption c are excluded from condition variable. This does not, however, correlate with the CE and KL terms of the loss which are lower for a condition-less setting. We attribute this to the model being transformed from a CVAE to a VAE, hence lifting the burden of capturing the conditional posterior distribution (i.e. the KL is now between an unconditional $q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x})$ and $\mathcal{N}(0, 1)$). Interestingly, however, excluding either the image or the caption achieves similar performance to when both are included, indicating that the caption acts as a good textual proxy of the image (a reassurance of our $sim_{c,q}$ metric).

C. Extended Quantitative Results on 2VD task

Extending Tab. 4 in the main paper, Tab. 7 here shows results for $\mathbf{B}/\mathbf{B}_{A\mathbf{R}}$ trained with permutations of the image *i* and caption *c*) (denoted by + if included in the condition, and – otherwise). We note the decrease in CE and KL as conditions (*i*, *c*) are excluded from the model. This is expected since the task of dialogue generation is made simpler without the constrains of an explicit visual/textual condition.

D. Network architectures and training

The following section provides detailed descriptions of the architectures of our models A, B and B_{AR} . The descriptions are dense but thorough. We also include further details of our training procedure. Where not explicitly noted, each convolutional layer is proceeded by a batch normalisation

Table 6: 1VD evaluation of B/B_{AR} on *VisDial* (v0.9) test set. Results show ranking of answer candidates based on the S_{w2v} scoring function. Note that d-qa indicates the iterative evaluation method when the *ground-truth* dialogue history is provided, while $d-q\hat{a}$, the iterative evaluation method when the ground-truth dialogue history is provided (see §4.2). The + and – indicate models *trained* with and without respective conditions, image *i* and caption *c*.

Method	i	c		CE	KLD	MR	MRR	R@1	R@5	R@10	$sim_{cap,q}$	sim
В	+	+	d– qa	18.87	4.36	28.45	0.2927	23.50	29.11	42.29	0.4374	2.68
	+	+		25.10	4.02	30.57	0.2188	16.06	20.88	35.37	0.4118	2.42
	_	+	d aâ	16.80	3.13	27.76	0.3243	26.59	33.21	47.65	0.4491	4.48
	+	_	a– qa	21.02	4.71	29.82	0.2144	15.25	21.07	34.96	0.4551	5.44
	_	_		19.35	13.34	29.00	0.3026	24.36	30.70	47.62	0.4638	6.17
B _{AR} 8	+	+	d– qa	15.11	2.53	25.87	0.3553	29.40	36.79	51.19	0.4703	4.30
	+	+		25.70	2.21	29.10	0.2864	22.52	29.01	48.43	0.3885	3.47
	_	+	d a â	16.19	2.80	26.04	0.3566	29.62	36.75	50.62	0.4626	4.17
	+	_	a– qa	20.39	2.89	28.99	0.3024	24.33	30.74	47.17	0.4461	8.16
	_	_		20.92	2.84	28.79	0.3045	24.46	30.99	48.10	0.4442	0.18
	+	+	d– qa	16.04	1.89	26.30	0.3422	28.00	35.34	50.54	0.4708	4.84
D 10	+	+		24.77	1.81	29.15	0.2869	22.68	28.97	46.98	0.4058	2.85
$\mathbf{B}_{\mathbf{AR}}$ 10	_	+	d a â	19.97	2.58	26.84	0.3212	25.90	32.92	47.68	0.4424	5.95
	+	_	a– qa	20.39	2.79	27.27	0.3157	25.45	32.26	47.87	0.4707	13.22
	-	-		19.17	0.00	29.00	0.3026	24.36	30.70	47.62	0.4614	0.00

layer (with momentum = 0.001 and learnable parameters) and a *ReLU* activation.

Prior network The prior neural network, parametrised by θ , takes as input the image *i*, the caption *c* and the dialogue context. For the model A, this context is h_t^+ , containing the dialogue history up to t-1 and the current question q_t . For models **B**/**B**_{AR}, the dialogue context is the null set $(h = \emptyset)$. To obtain the image representation, we scale and centre-cropped each image to $3 \times 224 \times 224$ and feed it through VGG-16 [23]. The output of the penultimate layer is extracted and ℓ_2 -normalised (as in [6]) to obtain a 4096dimensional image feature vector. For the caption, we pass cthrough a pre-trained word2vec [18] model (we do not learn these word embeddings) to obtain $\mathbf{\mathring{c}} \in \mathbb{R}^{300 \times L}$ where L is the maximum sentence length (L = 64). For the dialogue context (relevant only in the case of A) we pass the one-hot encoding of each word through a learnable word embedding module. We stack these embeddings as described in §3.1 of the main paper to obtain $\mathbf{h}_{t}^{+} \in \mathbb{R}^{E \times L \times K}$, where E is the word embedding dimension (E = 256), L is the maximum sentence length (L = 64) and K is the number of dialogue entries at time t. We encode these inputs convolutionally to obtain y (the encoded condition) as follows: \mathring{c} is passed through a convolutional block (output size $64 \times 8 \times 8$) and concatenated with the image feature vector (reshaped to $64 \times 8 \times 8$). The concatenated output is passed through a convolutional block to obtain the jointly encoded image-caption (output size $64 \times 8 \times 8$). If $h \neq \emptyset$, then the context is passed

through a convolutional block (output size $64 \times 8 \times 8$) and is concatenated with the encoded image-caption and passed through yet another convolutional block to get the encoded image-caption-context (output size $64 \times 8 \times 8$). We call this the encoded condition \boldsymbol{y} . The encoded condition \boldsymbol{y} is then passed through a further convolutional block (output size $256 \times 4 \times 4$) followed by two final convolutional layers (in parallel) to obtain $\boldsymbol{\mu}_p$ and $\log \sigma_p^2$, respectively, the parameters of the conditional prior $p_{\theta}(\boldsymbol{z} \mid \boldsymbol{y})$. At this stage, $\boldsymbol{\mu}_p$ and $\log \sigma_p^2$ are both of size $512 \times 1 \times 1$ (the latent dimensionality). At test time, a sample is obtained via $\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{z}; \boldsymbol{\mu}_p, \boldsymbol{\sigma}_p^2)$ and is passed to the decoder in order to generate a sample $\hat{\boldsymbol{a}}_t$ (for **A**) or $\hat{\boldsymbol{d}}$ (for **B/B_{AR}**).

Encoder network The encoder network, parametrised by ϕ , takes x as input along with the encoded condition, y, obtained from the prior network. For model A, x = a_t and $y_{-} = \{i, c, h_t^+\}$. For models B/B_{AR} , x = d = $\langle (\boldsymbol{q}_t, \boldsymbol{a}_t) \rangle_{t=1}^T$ and $\boldsymbol{y} = \{\boldsymbol{i}, \boldsymbol{c}\}$. In all models, \boldsymbol{x} is passed through a learnable word embedding module, and the word embeddings stacked (see $\S3.1$ in the main paper) to obtain $\mathring{\boldsymbol{x}} \in \mathbb{R}^{E \times L \times M}$, where E = 256, L = 64 and M is the number of entries in x (for A, M = 1 and for $B/B_{AR} M =$ 2T). In this way, we transform x into a single-channel answer 'image' in the case of A, and a multi-channel image of alternating questions and answers in the case of B/B_{AR} . \dot{x} is then passed through a convolutional block (output size $64 \times 8 \times 8$), the output of which is concatenated with y and forwarded through another convolutional block (output size

Table 7: 2VD evaluation on *VisDial* (v0.9) test set for $\mathbf{B}/\mathbf{B}_{AR}$ models. Note that d left blank indicates the block evaluation method, when a whole dialogue is generated given only an image and its caption, while $d-\hat{q}\hat{a}$ indicates the iterative evaluation method when previously generated questions *and* answers are included in the dialogue history (see Section 4.2). The + and – indicate models *trained* with and without respective conditions, image *i* and caption *c*.

Method	i	c	d	CE	KLD	$sim_{oldsymbol{c},q}$	sim
			Ø	31.18	4.34	0.4931	14.20
	+	+	d – $\hat{q}\hat{a}$	25.40	4.01	0.4091	1.86
			Ø	29.09	3.26	0.4889	11.23
В	-	+	d – $\hat{q}\hat{a}$	24.59	3.05	0.3877	3.45
	_		Ø	28.60	4.26	0.4634	15.56
	+	-	d – $\hat{q}\hat{a}$	29.85	4.66	0.4221	3.54
			Ø	19.92	6.42	0.4590	6.34
	-	-	d – $\hat{q}\hat{a}$	19.34	0.00	0.4638	0.00
			Ø	28.81	2.54	0.4878	31.50
	Ŧ	+	d – $\hat{q}\hat{a}$	26.60	2.29	0.3884	2.39
			Ø	30.59	2.72	0.4889	43.17
B _{AR} 8	-	Ŧ	d – $\hat{q}\hat{a}$	26.15	2.77	0.3758	3.57
	+		Ø	31.51	2.91	0.4602	24.75
	Τ		d – $\hat{q}\hat{a}$	21.41	2.68	0.4453	5.49
	_	_	Ø	20.32	2.77	0.4464	0.26
			d – $\hat{q}\hat{a}$	21.53	2.99	0.4419	0.10
	-	-	Ø	28.49	1.89	0.4927	44.34
	т 	т	d – $\hat{q}\hat{a}$	24.93	1.80	0.4101	2.35
	_	+	Ø	30.83	2.53	0.4951	38.60
B _{AR} 10		I	d – $\hat{q}\hat{a}$	28.59	2.52	0.3903	1.91
	+		Ø	30.18	2.89	0.4592	100.81
	Τ		d – $\hat{q}\hat{a}$	28.32	2.44	0.4334	6.73
	_	_	Ø	19.60	0.00	0.4585	0.00
			d – $\hat{q}\hat{a}$	19.17	0.00	0.4614	0.00

 $256 \times 4 \times 4$). This output is forwarded through two final convolutional layers (in parallel) to obtain μ_q and $\log \sigma_q^2$, the parameters of the conditional latent posterior $q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x}, \boldsymbol{y})$. Here μ_q and $\log \sigma_q^2$ are both of size $512 \times 1 \times 1$.

At train time, the KL divergence term of the ELBO is computed using $\{\mu_q, \sigma_q\}$ (from the encoder network) and $\{\mu_p, \sigma_p\}$ (from the prior network).

Decoder network The decoder network (for simplicity, the parameters of the prior and decoder network are subsumed into θ) takes as input a latent z and the encoded condition y. During training, z is sampled from a Gaussian parametrised by the μ_q and exponentiated $\log \sigma_q^2$ outputs of the *encoder* network. This distribution is $q_{\phi}(z \mid x, y)$. At test time, z is sampled from a Gaussian parametrised by the μ_p and exponentiated $\log \sigma_p^2$ outputs of the *prior* network. This distribution is $p_{\theta}(z \mid y)$. At both train and test time, we employ the commonly-used 're-parametrisation trick' [17] to compute the latent sample as $z = \mu + \epsilon \sigma$ where $\epsilon \sim \mathcal{N}(0, 1)$ and μ and σ correspond to those derived from the encoder or prior network as described above.

The sample z is then transformed through a transposeconvolutional block (output size $64 \times 8 \times 8$), concatenated with y and forwarded through a convolutional block (output size $64 \times 8 \times 8$). This output is forwarded through a second transpose-convolutional block, producing an intermediate output volume of dimension $M \times E \times L$ which we permute to match the size of \mathring{x} . As before, E = 256, L = 64 and M = 1 (for A) or M = 2T (for B/B_{AR}).

Following this, our models diverge in architecture: **A** and **B** employ a standard linear layer which projects the *E* dimension of the intermediate output to the vocabulary size *V*. The **B**_{AR} model instead employs an autoregressive module (detailed below) followed by this standard linear layer. At train time, the *V*-dimensional network output is *softmax*ed and used in the computation of the CE term of the ELBO. At test time, the *argmax* of the (*softmax*-ed) output is taken to be the index of the word token predicted. We share the weight matrices of the decoder's final linear layer and the encoder and prior's learnable word embedding module (which are the same size by virtue of our network architecture) with the motivation that language encoders and decoders should share common word representations.

Autoregressive block The autoregressive (AR) block (AR-N in Fig. 4 - bottom) in **B**_{AR}'s decoder is inspired by *PixelCNN* [26] which sequentially predicts the pixels in an image along the two spatial dimensions. In the same fashion, we use an autoregressive approach to sequentially predict the next sentence (question or answer) in a dialogue. Since our framework is convolutional with sentences viewable as 'images', our approach can similarly be adapted from that of [26, 9]. We first reshape the intermediate output of the decoder to $E \times L * M$ (essentially 'unravelling' the dialogue sequentially into a stack of its word embeddings). We then apply a size-preserving masked convolution to the reshaped output (followed by a learnable batch normalisation and a *ReLU* activation). We call this triplet an AR layer. The masked convolution of the AR layer ensures that future rows (i.e. future *E*-dimensional word embedding) are hidden in the prediction of the current row/word embedding. We apply N AR layers in this way with each layer taking in the output of the previous AR layer. Following the AR-N block, a linear layer projects the final output's E dimension to the vocabulary size V. We report numbers for $N = \{8, 10\}$. We base our implementation of the AR block on a publiclyavailable implementation of PixelCNN.

E. Dialogue preprocessing

The word vocabulary is constructed from the *VisDial* v0.9 [6] training dialogues (not including the candidate an-

swers). The dialogues are preprocessed as follows: apostrophes are removed, numbers are converted to their worded equivalents, and all exchanges are made lower-case and either padded or truncated to a maximum sequence length (L = 64). The vocabulary is also filtered such that words with a frequency of <5 are removed and replaced with the UNK token. After pre-processing and filtering, the vocabulary size is V = 9710.

F. Extended Qualitative Results

We present additional qualitative results for the **A** model in Figs. 4 and 5 (1VD task) and for the **B**_{AR}10 model (under the block evaluation setting) in Figs. 6 and 7 (2VD task). Note that for both, different colours indicate generations (\hat{a}_t for **A** and \hat{d} for **B**/**B**_{AR}) from different samples of *z*. In Figs. 6 and 7, whole generated dialogue blocks are shown with coloured sections indicating subsets exhibiting coherent question-answering and white sections indicating subsets that are not entirely coherent.

Question	Ground- truth answer	Z1	Z ₂	Z ₃	Z4	z ₅
How old is the girl?	Maybe three	I can't tell	Looks about six	She is maybe	I can't tell	Yes
What race is the girl?	White	All	Yes	White	White	Caucasian
Is she outside?	Yes	Yes	No	Yes	Yes	Yes
Is her hair long or short?	Short	Short	Short	Yes	Short	Short
What color is her hair?	Blonde	Blonde	Blonde	Blonde	Brown	Brown
Is her hair curly or straight?	It's straight	Straight	Straight	Short	Straight	Straight
What is she wearing?	Pink shirt, white pants	Clothes, shirt and top	T-shirt and jacket and pants	T-shirt and pants	Shirt and pants	Jeans like a pajamas
Is the teddy bear in her lap?	Yes	There	No	Yes	Yes	Yes
What color is the teddy bear?	White	Brown	Brown	Tan	Tan	Tan and white
Is it nice outside?	Yes sunny	Yes	Yes	It looks sunny	It looks lovely	Yes
Question	Ground- truth answer	z1	Z 2	Z 3	Z4	Z 5
How old does she look?	Around seven or eight	I cannot tell about her	Looks about six	She i: abou teenage aroun	s t I can r tei d	't I can' ll te
Any other people?	No	No	Yes	Ye:	s I	No I don see anyon
Any buildings?	No	Yes	No	N	o Secre	et Ye
Is it day or night?	Day	It looks like it in image is in so	Daytime	Da:	y Da	ay Ye
Is it raining?	No	Yes	Yes	N	0 1	No I
What color umbrella?	Pink and clear	Dark colored	White	White	e Gre	ey Whit

A young girl holding an umbrella on the sidewalk

Question	Ground- truth answer	Z1	Z2	Z3	Z4	2 ₅
Are there any animals present?	No animals	No	No	Yes	No	No
Are there any kids?	Two kids	Yes	No	Yes	Yes	I
Boys or girls?	Boys	Воу	Boys	I	Boys	Yes
Are there any adults?	One adult	Yes	Yes	Yes	Yes	Yes
What are they doing?	Adult is flying kite and boys are watching	Playing	I can't see their faces	She is watching parasols	Spectating	They are playing watching
Is it sunny?	Yes	Yes	No	I can't tell	Yes	Yes
Are there trees?	Yes	Yes	Yes	Yes	Yes	There are no background
Is kite colorful?	Yes	No	Yes	Yes	Yes	Yes
Is kite in air?	Yes	Yes	Yes	No	Yes	No
Are boys smiling?	Yes	Yes	Yes	No	Yes	Yes

Open

Yes

Brown

long

I'd say

It is open

No

Brown

Yes

Open

Yes

Brown

Short

Is it open or closed?

Is it sunny?

What color is her hair?

Is it long or short?

Oper

Long

I can't tell

Dark brown

Yes

Yes

Brown

Short

Open

Brown

Short

No

Figure 4: Examples of diverse answer generations from the **A** model for the 1VD task.

Is the man in a park?	I can te	't It appe	so	Yes	Yes	Yes	He is in	
How old is the man?	Thirt	ys Loo	cate Maybe	late	Thirtys	Thirty	Twentys	. whiteen
Has the frisbee left his hand?	Ye Ye	es	Yes	Yes	Yes	Yes	Yes	
Can you tell what he's throwing it towards?	1	No	No tel	can't 1, is ball	I can't tell	No	Yes	
What does 'casually thrown' mean?	It loo like l just toss	ks he ed	At ally cas	At least relief casually		He has in a kateboard	Casually	
What color is the frishee?	Whi	it te Wi	hite	White	Blue	shirt It is white	White	
What is the man wearing?	A whit scree	te en T-si	hirt T- and s	shirt H horts	ie looks casual	Shorts	T-shirt and shorts	
Is he wearing jeans?	I can't se his pan	ee ts	No S	horts	Yes	No	Yes	
Are there other people around?	1	No	No	No	No	No	No	
Do you see any animals?	, 1	No	No	No	No	No	No	A man casually throws a frisbee into the air
Question	Ground- truth answer	\mathbf{z}_1	Z ₂	z,	Z4	z 5		
Is it a tall wall?	Yes, it goes	Yes	s Yes	Yes can'	t Ne	D N	0	
Is there only one polar bear?	Yes, just	Yes	s One	Ye	s One	e Ye	s	-
Any people there?	No	yes	I can't see any	N	o Ne	o Ye	s	
Do you see quite a lot of water?	Yes the polar bear is underwater	Yes	I can't see the ground	N	o Yes can' see an	t Yes can' y tel	t 1	and the second sec
Can you see the sky?	No	n No	o Yes	Ye	s Ye:	s Ye	s	
Is this picture in color?	Yes it is	Yes	s Yes	Ye	s Ye:	s Ye	s	and the second s
Is this a professional picture?	I don't think it is	Looks like	e I	Ye	s N	o Ye	s	R ASIN HALLES
Do you see any buildings?	No	n No	o Yes	There ar n building	e o :	I N	•	
Any trees?	No	No	There see	Ye	s Ye:	s N	•	
Is this a close- up type of photo?	No	Yes	s Yes	Ye	s Ye:	It look s like it i in fram	s e	A polar bear swimming in water near a rock wall
Question	Ground- truth answer	z,	Z ₂	\mathbf{z}_{3}	Z4	z 5		
How many children are there?	Many	There are six	One	Three	e Ther	e More that ten te	an en	
All boys or girls?	Both	I can't tell	Воу	Girl an eac	i It appear boy gir	s They and	ce so	
How old do they look?	Six	Five	Ten-twelve	Twent	7 Ye	s Five lool s:	cs Lx	
Are they outside?	Yes	I	Yes	Ye	з Уе	s Ye	25	
Is it sunny out?	Maybe	Yes	Yes	Ye	a Ye	s Ye	25	
Are they in a field?	Yes	Yes	No	N	y Ye	s 1	10	
Are there many adults around?	Zero	No	Yes	N	Ye Ye	s 1	10	
Is there only one ball?	Yes	Yes	Yes	Ye	a Ye	s Ye	25	
Are there soccer nets?	Not visible	No	No	Ye	I can' sa	t y	es	
What's in the background?	Fence	Trees	I see the grass	There are numerou snow trees	a Lot	s I see sor gras	ne ss Sma	all children in red and blue uniforms, kicking a red soccer ball

Z4

 \mathbf{z}_{5}

 \mathbf{Z}_2

 \mathbf{z}_3

 \mathbf{Z}_1

Groundtruth answer

Question

Figure 5: Examples of diverse answer generations from the A model for the 1VD task – continued.

		How many sh there?	neep are	Four	Is this a color?		Yes	Can you see any people?	No
		Are you any	y any water	? No	Can you see any?	1	No	Is the photo in color?	Yes
		What color	are they?	Brown	Can you see trees	s? 1	No	What is the fence made?	I is
		Any people?	?	No	Any people?	1	No	Can you see any?	No
		What are th	ne the?	Grazi	ng Is it sunny?	I	No	How many sheep are there?	I are
		Are there p	people?	No	Are it sunny?		Yes	What color is the shacks?	It
	B Alar Gito Guo D Gito Gito Alar	Is it sunny	/?	Yes	Any people?	I	No	Does the grass have?	Yes
		Can any ani	imals?	No	Any people?	1	No	Can you see any sky?	No
	o de la serie	Are any any	/?	No	What time the day	/?	Can't tell	Is the grass green?	Yes
Sheep standing near orange netting in grassy fie	ald	Is any sunr	ıy?	Yes	Are they eating?		Yes	Do you see any sky?	No
		Is the people m	ale?	(es	Can you see the bal	1? Ye	s, I	Is this a professional game?	Yes
		Is the person w	earing? N	(es	Is the player a?	Yei	s, is	Any many players are?	One
		Can you see any	sky? M	10	Is it a a a?	Yei	s, is	Is the a?	Yes
	d	What color is t uniforms?	he	Vhite €	Can you see the baseball?	Yei	s, is	Can you see the pitcher?	No
A Contraction of the second se		Is there grass	in? Y	(es	What is the the wearing?	It we	is aring	Is the a in?	Yes
		How many people there?	are	Ewo	How is the the bat?	I	is is	Is it a in the?	Yes
		Is it sunny sun	ny? Y	(es	Is there any other people?	Ye	s	What color is the batters?	White
		Can you see any	? 1	(es	Can you see the pitcher?	No	, the e	Is there people in the?	No
	What color is th		The	Can you see the scoreboard?	No		Is the pitcher in?	Yes	
A baseball player for the Chicago Cubs stands at hom	me plate	Is it sunny?	2 Z	íes 🛛	Can you see the sky	? No		Is it sunny?	Yes
	What color is?	N	White	What color is the fridge?		ite	Can you see the?	No	
		Is there a fridg		'es	What color is?	Wh	ite	What appliances?	Microwave
		Is there any people Is it a? Is there a window Is there a? Is there a? Is the fridge on		10	Can you see the sin	k? Ye:	s	Is the fridge on?	Yes
SH OF	-			'es	Is there a window?	No		Is it a?	Yes
				10	Any people?	No		Is the fridge have magnets?	Yes
				10	Is there?	No		Can you see what time?	Colgate
				'es	How pics the?	No		Is there a on the?	No
				'es	What color is the walls?	Wh:	ite	Can you tell what time?	No
		Can you see the	window? N	10	Any windows?	No		Is there a?	Yes
An image of a kitchen loft style setting		Is it a?	У	'es	How about on?	Two	0	Are there any people?	No
	Is the pho	oto in color?	Yes	What c unifor	olor are the ms?	One is white	Wł	hat is the the of?	One
	Is it a pr photo?	rofessional	Yes	Are th	ey wearing?	Yes	A	e there see other	Yes
	Are the me	en wearing or?	Yes	Are th	ere a?	Yes	A1	they people	Yes
	Is any of	wearing?	Yes	Is the	re a?	Yes	Is	they wearing the they wearing the they wearing the they wear the	Yes
	Are there the?	any people in	Yes	Is the	re trees?	Yes	Whur	hat color are the hiforms?	White
	How many r	nen?	Two		sunny?	Yes	Do	es the have wearing	Yes
	Are the me wearing un	en wearing hiforms?	Yes	Are th wearin	e wearing g?	Yes they	Wł we	at color the aring?	A
The second s	Are there people?	see other	No	Are th cleats	ey wearing ?	Yes	Ca	an you see any ball?	No
Two guys playing baseball.	What color uniforms?	r are the	White	Are th	ere a fence?	Yes	Ca	an you see any baches?	No
with trees in the back	Do the hav	ve the the?	Yes	Can yo	u see the?	No	Do	you see the color?	Yes

Figure 6: Diverse two-way dialogue generations from the B_{AR} 10 model (block evaluation) for the 2VD task.

State State	Is the pic color?	ture in	Yes	How ma there?	ny people are	Two	W m	hat color is the an's?	White								
	Is the pho	to in?	Yes	What a	re they?	I	W	hat is the man man?	I								
	What color frisbee?	What color is the frisbee? Can there?		What color is the frisbee?		color is the		What color is the frishee?		hat color is the risbee?		Can you see the frisbee?		Yes	н	ow old is the man?	I
	Can there?			What c frisbe	olor of the e?	Blue	с	an you see any ball?	Yes								
	Is the gra	ss visible?	Yes	Is the	re people?	Yes	I	s there any other eople?	No								
	Is the sun	ny?	Yes	Is it	sunny?	Yes	D	oes the man have a?	No								
	Is it sunn	y?	Yes	What c	olor is the?	I	I	s the man wearing a at?	No								
	Is there o	ther people?	No	Is the	re sunny?	Yes	D	oes he man have a?	No								
A AL	What color frisbee?	is the	White	How ma there?	ny people are	ire I		re there see any ther people?	No								
Man and a boy playing ball in the grass	Is the man	Yes	Is the	y having fun?	Yes	с	an you see any sky?	No									
		Is it big?		Yes Are there only one?		2 1	ſes	What kind of dog	It looks								
		How is of the i	is?	I	What kind of dog is		I is like	Does he have a	Yes								
		Can you see the	e breed?	No	Does it look like b	be? I	It, it	Collar? Can you see what breed?	No								
	A	Are there any p	people?	No	Is there any people the boat?	e in N	No	Is it people boat the the boat?	No								
E-820 CR MS 815 Δ.)		Is the dog on?		No	What time of day is it? It		It is light	Is there oars?	No								
		Can you tell th of it?	ne breed	No	What color is the b	boat? W	White	How the boat the the?	I								
	9-	What is the col		It is	How it is like is?	1	I is	Is it sunny?	Yes								
		driver?	ine	No Is it a or or? Yes What is of day is?		1	It Is	How many barges?	No								
		Is it a Harley?	?			1	It is	Any people?	No								
A dog sitting on the inside of a w	hite boat	Any other?		No	What color is the b	boat?	It	Any?	No								
	What color	is the man?	White	Is thi	s a color photo?	Yes	W	hat is the color?	Is								
	What color	is his hair?	Black	What color is the tennis?		It	I	s the player?	Yes								
	How old is	he?	Thirtys	Are th people	Are there any other people in the?		I	s it a?	Yes								
	Is he have	other?	Yes	Can yo	u see any tennis?	No	I	s it outdoors?	Yes								
	Is it sunn	Υ?	Yes	Can yo in the	u see any people ?	No	I	s there people?	No								
A A A A A A A A A A A A A A A A A A A	Can you se	e any?	Yes	What t	ime of day is it?	It	I	s he playing?	Yes								
	Is it cour	Yes	Is the be?	appear like be	It	I	s it a?	Yes									
	Is it sunn	y?	Yes	Are th the th	ere any other in e?	No	I	s it a?	Yes								
	Is it sunn	y?	Yes	Can yo	u see the time?	No	с	an you see the net?	No								
A person is holding a ball and tennis racket	Can you see sky?		No	Can you see the season?		No	I	s it man?	Yes								

Figure 7: Diverse two-way dialogue generations from the B_{AR} 10 model (block evaluations) for the 2VD task – continued