

Large-scale Pretraining for Visual Dialog: A Simple State-of-the-Art Baseline

Vishvak Murahari¹ Dhruv Batra^{1,2} Devi Parikh^{1,2} Abhishek Das^{1→2}

¹ Georgia Institute of Technology ² Facebook AI Research

{vishvak.murahari,dbatra,parikh,abhshkdz}@gatech.edu

Abstract. Prior work in visual dialog has focused on training deep neural models on VisDial in isolation. Instead, we present an approach to leverage pretraining on related vision-language datasets before transferring to visual dialog. We adapt the recently proposed ViLBERT model for multi-turn visually-grounded conversations. Our model is pretrained on the Conceptual Captions and Visual Question Answering datasets, and finetuned on VisDial. Our best single model outperforms prior published work by $> 1\%$ absolute on NDCG and MRR.

Next, we find that additional finetuning using “dense” annotations in VisDial leads to even higher NDCG – more than 10% over our base model – but hurts MRR – more than 17% below our base model! This highlights a trade-off between the two primary metrics – NDCG and MRR – which we find is due to dense annotations not correlating well with the original ground-truth answers to questions.

Keywords: Vision & Language, Visual Dialog

1 Introduction

Recent years have seen incredible progress in Visual Dialog [1–22], spurred in part by the initial efforts of Das *et al.* [2] in developing a concrete task definition – given an image, dialog history consisting of a sequence of question-answer pairs, and a follow-up question about the image, to predict a free-form natural language answer to the question – along with a large-scale dataset and evaluation metrics. The state-of-the-art on the task has improved by more than 20% absolute ($\sim 54\% \rightarrow \sim 74\%$ NDCG) and the original task has since been extended to challenging domains, *e.g.* video understanding [23], navigation assistants [24–26].

While this is promising, much of this progress has happened in isolation, wherein sophisticated neural architectures are trained and benchmarked solely on the VisDial dataset. This is limiting – since there is a significant amount of shared abstraction and visual grounding in related tasks in vision and language (*e.g.* captioning, visual question answering) that can benefit Visual Dialog – and wasteful – since it is expensive and dissatisfying to have to collect a large-scale dataset for

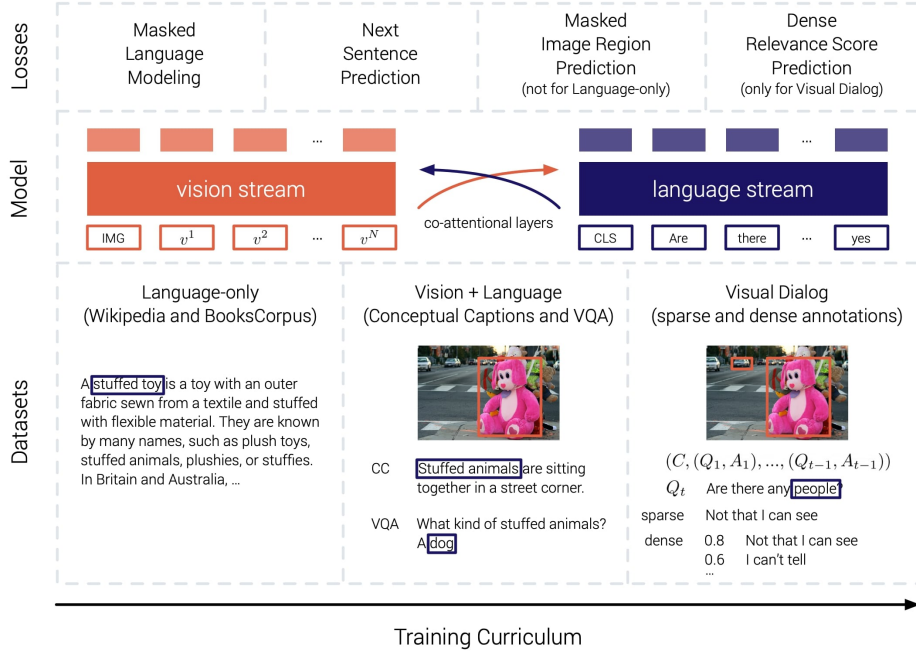


Fig. 1: First, the language stream of our model is pretrained on English Wikipedia and the BooksCorpus [27] datasets with the masked language modeling (MLM) and next sentence prediction (NSP) losses. Next, the entire model is trained on the Conceptual Captions [28] and VQA [29] datasets with the masked image region (MIR), MLM and NSP losses. Finally, we finetune the model on sparse annotations from VisDial [2] with the MIR, MLM and NSP losses, and optionally finetune on dense annotations.

every new task. In this work, we explore an approach to pretrain our model on other related vision and language datasets and then transfer to Visual Dialog.

Our work is inspired by prior work in transfer learning in computer vision and natural language understanding where large models [30–40] are pretrained on large datasets [27, 41, 42] with simple yet powerful self-supervised objectives to learn powerful representations that are then transferred to downstream tasks, leading to state-of-the-art results on a variety of benchmarks [41, 43]. Recent work has extended this to vision and language tasks [44–50], leading to compelling results in Visual Question Answering [29], Commonsense Reasoning [51], Natural Language Visual Reasoning [52], Entailment [53], Image-Text Retrieval [54, 55], Referring Expressions [56], and Vision-Language Navigation [57].

In this work, we adapt ViLBERT [44] to Visual Dialog. ViLBERT uses two Transformer-based [34] encoders, one for each of the two modalities – language and vision – and interaction between the two modalities is enabled by co-attention layers *i.e.* attention over inputs from one modality conditioned on inputs from the other. Note that adapting ViLBERT to Visual Dialog is not trivial. The Visual Dialog dataset has image-grounded conversation sequences that are

up to 10 rounds long. These are significantly longer than captions (which are ≤ 2 sentences) from the Conceptual Captions dataset [28] or question-answer pairs from VQA [29] used to pretrain ViLBERT, and thus requires a different input representation and careful reconsideration of the masked language modeling and next sentence prediction objectives used to train BERT [35] and ViLBERT [44].

This adapted model outperforms prior published work by $> 1\%$ absolute and achieves state-of-the-art on Visual Dialog. Next, we carefully analyse our model and find that additional finetuning on ‘dense’ annotations¹ *i.e.* relevance scores for all 100 answer options corresponding to each question on a subset of the training set, highlights an interesting trade-off – the model gets to $\sim 74.5\%$ NDCG (outperforming the 2019 VisDial Challenge winner), but an MRR of $\sim 52\%$ ($\sim 17\%$ below our base model!). We find this happens because dense annotations in VisDial do not correlate well with the ground-truth answers to questions, often rewarding the model for generic, uncertain responses.

Concretely, our contributions are as follows:

- We introduce an adaptation of the ViLBERT [44] model for Visual Dialog, thus making use of the large-scale Conceptual Captions [28] and Visual Question Answering (VQA) [29] datasets for pretraining and learning powerful visually-grounded representations before finetuning on VisDial [2]. Since captioning and VQA differ significantly from Visual Dialog in input size (≤ 2 sentence descriptions *vs.* ≤ 10 question-answer rounds), this requires rethinking the input representation to learn additional segment embeddings representing questions-answer pairs. Our adapted model improves over prior published work by $> 1\%$ and sets a new state-of-the-art.
- We next finetune our model on dense annotations *i.e.* relevance scores for all 100 answer options corresponding to each question on a subset of the training set, leading to even higher NDCG – more than 10% over our base model – but hurting MRR – more than 17% below our base model! This highlights a stark trade-off between the two primary metrics for this task – NDCG and MRR. Through qualitative and quantitative results, we show that this happens because dense annotations do not correlate well with the original ground-truth answers, often rewarding the model for generic, uncertain responses.
- Our PyTorch [58] code is publicly available² to encourage further work in large-scale transfer learning for VisDial.

2 Related Work

Our work is related to prior work in visual dialog [1–22], and self-supervised pretraining and transfer learning in computer vision and language [30–40].

¹ publicly available on visualdialog.org/data.

² github.com/vmurahari3/visdial-bert/

Visual Dialog. Das *et al.* [2] and de Vries *et al.* [1] introduced the task of Visual Dialog – given an image, dialog history consisting of a sequence of question-answer pairs, and a follow-up question, predict a free-form natural language answer to the question – along with a dataset, evaluation metrics, and baseline models. Follow-up works on visual dialog have explored the use of deep reinforcement learning [3, 4, 17], knowledge transfer from discriminative to generative decoders [5], conditional variational autoencoders [6], generative adversarial networks [7], attention mechanisms for visual coreference resolution [9, 11], and modeling the questioner’s theory of mind [10]. Crucially, all of these works train and evaluate on the VisDial dataset *in isolation*, without leveraging related visual grounding signals from other large-scale datasets in vision and language. We devise a unified model that can be pretrained on the Conceptual Captions [28] and VQA [29] datasets, and then transferred and finetuned on VisDial.

Self-Supervised Learning in Vision and Language. Building on the success of transfer learning in natural language understanding [33–40] leading to state-of-the-art results on a broad set of benchmarks [41, 43], recent work has extended this to vision and language tasks [44–50]. These works pretrain single [45, 48, 49] or two [44, 46]-stream Transformer [34]-based models with self-supervised objectives, such as next-sentence prediction and masked language/image modeling, on large-scale image-text datasets and have led to compelling results in Visual Question Answering [29], Commonsense Reasoning [51], Natural Language Visual Reasoning [52], Entailment [53], Image-Text Retrieval [54, 55], and Referring Expressions [56], and Vision-Language Navigation [57].

3 Adapting ViLBERT [44] for Visual Dialog

Lu *et al.* [44] introduced ViLBERT³, which extended BERT [35] to a two-stream multi-modal architecture for jointly modeling visual and linguistic inputs. Interaction between the two modalities was enabled through co-attention layers, *i.e.* attending to one modality conditioned on the other – attention over language conditioned on visual input, and attention over image regions conditioned on linguistic input. This was operationalized as swapping the key and value matrices between the visual and linguistic Transformer [34] blocks. We next discuss our changes to adapt it for Visual Dialog followed by our training pipeline.

Input Representation. Recall that the model gets image I , dialog history (including image caption C) $H = (C, (Q_1, A_1), \dots, (Q_{t-1}, A_{t-1}))$, question Q_t , and a list of 100 answer options $A_t = \{A_t^{(1)}, A_t^{(2)}, \dots, A_t^{(100)}\}$ as input, and is asked to return a sorting of A_t . We concatenate the t rounds of dialog history and follow-up question Q_t , with each question and answer separated by a <SEP> token. The overall input to the language stream is represented as:

$$\langle \text{CLS} \rangle C \langle \text{SEP} \rangle Q_1 \langle \text{SEP} \rangle A_1 \langle \text{SEP} \rangle, \dots, \langle \text{SEP} \rangle Q_t \langle \text{SEP} \rangle A_t \langle \text{SEP} \rangle \quad (1)$$

³ along with code released at github.com/jiasenlu/ViLBERT_beta.

Similar to Wolf *et al.* [59], we use different segment embeddings for questions and answers to help the model distinguish between the two and understand question and answer boundaries in the input. Captions and answers share the same segment embeddings. To represent the image, we follow [44, 60] and extract object bounding boxes and their visual features for top-36 detected objects in the image from a Faster R-CNN [61] (with a ResNet-101 [30] backbone) object detection network pretrained on the Visual Genome dataset [42]. The feature vector for each detected object is computed as mean-pooled convolutional features from the regions of that object. A 5-d feature vector, consisting of normalized top-left and bottom-right object coordinates, and the fraction of image area covered, is projected to the same dimensions as the feature vector for the detected object, and added to it, giving us the final visual features $\{v_1, \dots, v_{36}\}$. The beginning of this image region sequence (consisting of object detection features) is demarcated by an IMG token with mean-pooled features from the entire image. The overall input to ViLBERT can be written as the following sequence:

$$\langle \text{IMG} \rangle v_1, \dots, v_{36} \langle \text{CLS} \rangle C \langle \text{SEP} \rangle Q_1 \langle \text{SEP} \rangle A_1 \langle \text{SEP} \rangle, \dots, \langle \text{SEP} \rangle Q_t \langle \text{SEP} \rangle A_t \langle \text{SEP} \rangle \quad (2)$$

3.1 Pretraining on Conceptual Captions [28]

To pretrain the model, we follow [44] and train on the Conceptual Captions (CC) dataset, which is a large corpus (with $\sim 3\text{M}$ samples) of aligned image-caption pairs. During pretraining, the sum of the *masked language modeling* (MLM) loss [35] and the *masked image region* (MIR) loss is optimized. To compute the MLM loss, a set of tokens in the input sequence are masked and the model is trained to predict these tokens given context. We mask around 15% of the tokens in the input sequence. For the MIR loss, similar to the MLM loss, we zero out 15% of the image features and the model learns to predict the semantic category of the masked out object (out of 1601 classes from Visual Genome [42, 60]).

3.2 Pretraining on VQA [29]

The VQA dataset is quite related to Visual Dialog in that it can be interpreted as independent visually-grounded question-answer pairs with no dialog history, and thus is a natural choice for further pretraining prior to finetuning on VisDial. Similar to Lu *et al.* [44], we pretrain on VQA by learning a small decoder – a two-layer MLP – on top of the element-wise product between the image and text representations to predict a distribution over 3129 answers.

3.3 Finetuning on Visual Dialog [2]

To finetune on Visual Dialog, we use the MLM loss along with the next sentence prediction (NSP) and MIR losses. For MLM, we mask 10% of the tokens

in the dialog sequence. For MIR, similar to pretraining, we mask 15% of the image features. Note that the discriminative task in visual dialog is to identify the ground-truth answer from a list of 100 answer options consisting of popular, nearest neighbors, and random answers from the dataset. We achieve this through the NSP loss. The NSP head is trained to predict 1 when the ground-truth answer is appended to the input sequence, and 0 when a negative answer sampled from the remaining answer options is appended to it. Each image in VisDial has 10 rounds of dialog, leading to 10 sets of positive and negative samples for the NSP loss per mini-batch. Since these are fairly correlated samples, we randomly sub-sample 2 out of these 20 during training. At test time, we use log-probabilities from the NSP head to rank the 100 answer options per round.

3.4 Finetuning with Dense Annotations

The authors of [2] recently released dense annotations⁴ *i.e.* relevance scores for all 100 answer options from A_t corresponding to the question on a subset of the training set. These relevance scores range from 0 to 1 and are calculated as the ratio of number of human annotators who marked a particular answer option as correct to the total number of human annotators ($= 4$). So 1 means that the answer option was considered correct by 4 human annotators. In our final stage of training, we utilize these dense annotations to finetune our model. Concretely, we use the NSP head to predict likelihood scores $\hat{\ell}_t^{(i)}$ for each answer option $A_t^{(i)}$ at round t , normalize these to form a probability distribution over the 100 answers $\hat{y}_t = [\hat{y}_t^{(1)}, \dots, \hat{y}_t^{(100)}]$, and then compute a cross-entropy (CE) loss against the normalized ground-truth relevance scores y_t , given by $-\sum_i y_t^{(i)} \log \hat{y}_t^{(i)}$.

4 Experiments

To compare to previous research, we conduct experiments on VisDial v1.0 [2]. The dataset contains human-human dialogs on $\sim 130k$ COCO [62]-like images. We follow the original splits and use $\sim 120k$ for training, $\sim 2k$ for validation, and $\sim 8k$ for testing. We next describe the various settings we experiment with.

Evaluation Metrics. We use metrics introduced in [2]. Specifically, given the predicted ranking of 100 answer options from a model at each round, we compute retrieval metrics – mean rank (MR) of the ground-truth answer, mean reciprocal rank (MRR), and recall@ k ($k = \{1, 5, 10\}$). Additionally, along with the release of dense annotations, *i.e.* relevance scores $\in [0, 1]$ for all 100 answer options, a new metric – NDCG – was introduced. NDCG accounts for multiple correct answers in the option set and penalizes low-ranked but correct answer options.

⁴ publicly available on visualdialog.org/data.

4.1 Language-only

We begin with a ‘blind’ setting, where given the dialog history and follow-up question, and without access to the image, the model is tasked with predicting the answer. We do not use the ViLBERT formulation for these experiments, and finetune the BERT model released in [35] and pretrained on BooksCorpus [27] and English Wikipedia. For the MLM loss, we mask 15% of tokens and subsample 8 out of 20 sequences per mini-batch during training. We experiment with two variants – training only with NSP, and training with both NSP and MLM. See Table 3 for language-only results (marked ‘L-only’). This setting helps us benchmark gains coming from switching to Transformer [34]-based architectures before the added complexity of incorporating visual input.

Varying number of dialog rounds. We train ablations of our language-only model (with NSP and MLM losses) where we vary the number of rounds in dialog history, starting from 0, where the input sequence only contains the follow-up question and answer, to 2, 4, and 6 and 10 rounds of dialog history (Table 1).

Zero-shot and ‘cheap’ finetuning. We report performance for ablations of our NSP+MLM model with no/minimal training in Table 2. First, we do a zero-shot test where we initialize BERT with weights from Wikipedia and BooksCorpus pretraining and simply run inference on VisDial. Second, with the same initialization, we freeze all layers and finetune only the MLM and NSP loss heads.

4.2 Finetuning on VisDial

We finetune ViLBERT on VisDial with four different weight initializations – 1) with randomly initialized weights, 2) from the best language-only weights (from Section 4.1) for the language stream (visual stream and co-attention layers initialized randomly), 3) from a model pretrained on CC [28] (as described in Section 3.1) and 4) from a model pretrained on CC [28]+VQA [29] (as described in Section 3.2). 1) helps us benchmark improvements due to pretraining, 2) helps us benchmark performance if the model learns visual grounding solely from VisDial, 3) quantifies effects of learning visual grounding additionally from CC, and 4) helps us quantify improvements with additional exposure to visually-grounded question-answering data. See Table 3 for results.

4.3 Finetuning with Dense Annotations

Finally, we finetune our best model from Section 4.2 – marked ‘w/ CC+VQA’ in Table 3 – on dense annotations, as described in Section 3.4. Note that computing the CE loss requires a separate forward pass for each of the 100 answer options, since dialog history, question, answer are all concatenated together before passing as input. This is memory-expensive, and so in practice, we subsample and only use 80 options, and use gradient accumulation to (artificially)

construct a larger mini-batch. Finetuning with the CE loss only leads to significant improvements on NDCG but hurts other metrics (see Table 3). We discuss and analyse this in more detail later. But to control for this ‘metric-overfitting’, we also train a variant with both the CE and NSP losses.

5 Results

We list findings from all experiments described in Section 4 below.

Table 1: Performance of the NSP + MLM language-only model on VisDial v1.0 val as the number of dialog history rounds is varied

# history rounds	NDCG \uparrow	MRR \uparrow	R@1 \uparrow	R@5 \uparrow	R@10 \uparrow	MR \downarrow
0	50.54	54.29	38.88	72.67	83.09	5.90
2	53.69	61.31	46.83	78.96	88.15	4.51
4	55.10	62.83	48.36	80.61	89.57	4.19
6	55.69	63.73	49.31	81.13	90.06	4.04
10	57.22	64.10	50.05	81.09	90.00	4.16

Table 2: Performance of the NSP + MLM language-only model on VisDial v1.0 val with no / minimal training (described in Sec. 4.1)

Model	NDCG \uparrow	MRR \uparrow	R@1 \uparrow	R@5 \uparrow	R@10 \uparrow	MR \downarrow
No training	11.63	6.88	2.63	7.17	11.30	46.90
Loss heads only	19.69	9.81	3.42	10.44	18.85	31.38

- **Language-only performs well.** The language-only model gets to 57.22 on NDCG and 64.10 on MRR (Table 3), which is already competitive with several prior published works (Table 4). These trends are consistent with high human performance on VisDial [2] with just language (question and dialog history) – 48.5 on MRR – which further improves to 63.5 on MRR with image.
- **Increasing dialog history rounds helps.** We report performance of the language-only model as a function of dialog history rounds in Table 1 and Fig. 2a. Note that the change in performance from including 0 to 4 rounds of dialog history (+4.56 on NDCG, +8.54 on MRR) is much more than from 4 to 10 dialog history rounds (+2.12 on NDCG, +1.27 on MRR). Thus, performance continues to go up with increasing dialog history rounds but starts to plateau with ≥ 4 history rounds. We believe these improvements are largely indicative of the Transformer’s ability to model long-term dependencies.
- **Zero-shot model performs poorly.** Running inference with the language-only model pretrained on BooksCorpus [27] and Wikipedia without any finetuning on VisDial only gets to 11.63 on NDCG and 6.88 on MRR (Table 2). Finetuning the loss heads with all other layers frozen leads to an improvement

- of ~ 8 NDCG points over this. This low performance can be attributed to significantly longer sequences in VisDial than the model was pretrained with.
- **VQA initialization helps more than random or CC initialization.** Finetuning ViLBERT on VisDial with weights initialized from VQA pretraining gets to 64.82 on NDCG and 68.97 on MRR, ~ 3 points better than random initialization on NDCG and ~ 2 points better than CC pretraining (Table 3). We believe poorer transfer from CC is because both VQA and VisDial have images from COCO and are more closely related tasks than captioning on CC.
 - **Dense annotations boost NDCG, hurt MRR.** Finetuning with the CE loss leads to 74.47 on NDCG – a $\sim 10\%$ improvement over the ‘w/ CC + VQA’ base model – but 50.74 on MRR, a $\sim 17\%$ decline below the base model (Table 4). This is a surprising finding! We carefully analyze this behavior in Section 6.
 - **Ensembling does not improve performance.** We trained 3 models initialized with different random seeds for each of the 3 variants (‘w/ CC + VQA’, ‘CE’ and ‘CE + NSP’) and aggregated results by averaging the normalized scores from the 3 models. We did not observe any significant improvement.

Table 3: Results on VisDial v1.0 val (with 95% CI). \uparrow indicates higher is better.

	Model	NDCG \uparrow	MRR \uparrow	R@1 \uparrow	R@5 \uparrow	R@10 \uparrow	MR \downarrow
L-only	NSP	55.80 ± 0.9	63.37 ± 0.5	49.28 ± 0.7	80.51 ± 0.5	89.22 ± 0.4	4.32 ± 0.1
	NSP + MLM	57.22 ± 0.9	64.10 ± 0.5	50.05 ± 0.7	81.09 ± 0.5	90.00 ± 0.4	4.16 ± 0.1
+vision	Random init	61.88 ± 0.9	67.04 ± 0.5	53.51 ± 0.7	83.94 ± 0.5	92.27 ± 0.4	3.55 ± 0.1
	w/ L-only	62.08 ± 0.9	67.73 ± 0.5	54.67 ± 0.7	84.02 ± 0.5	92.07 ± 0.4	3.58 ± 0.1
	w/ CC [28]	62.99 ± 0.9	68.64 ± 0.5	55.55 ± 0.7	85.04 ± 0.5	92.98 ± 0.4	3.36 ± 0.1
	w/ CC [28] + VQA [29]	64.82 ± 0.9	68.97 ± 0.5	55.78 ± 0.7	85.34 ± 0.5	93.11 ± 0.4	3.35 ± 0.1
+dense	CE	75.10 ± 1.1	52.12 ± 0.6	39.84 ± 0.7	64.93 ± 0.7	80.47 ± 0.5	6.26 ± 0.1
	CE + NSP	69.11 ± 1.0	65.76 ± 0.5	53.30 ± 0.7	80.77 ± 0.5	90.00 ± 0.4	4.33 ± 0.1

We report results from the Visual Dialog evaluation server⁵ for our best models – ‘w/ CC + VQA’, ‘CE’ and ‘CE + NSP’ – on the unseen `test-std` split in Table 4. We compare against prior published results and top entries from the leaderboard. Our models outperform prior results and set a new state-of-the-art – ViLBERT with CC + VQA pretraining on MRR, R@k, MR metrics, and further finetuning with a CE loss on dense annotations on NDCG. Finally, adding NSP loss along with CE (as in Section 4.3) offers a balance between optimizing metrics that reward both sparse (original ground-truth answers) and dense annotations.

6 Analysis

As described in Section 5, finetuning on dense annotations leads to a significant increase in NDCG, but hurts the other 5 metrics – MRR, R@1, R@5, R@10 and

⁵ evalai.cloudcv.org/web/challenges/challenge-page/161/leaderboard/483

Table 4: Results on VisDial v1.0 test-std. \uparrow indicates higher is better. \downarrow indicates lower is better. \dagger denotes ensembles. Best single-model results are **bolded** and best ensemble results are underlined. \star denotes the winning team of the 2019 Visual Dialog Challenge.

	Model	NDCG	\uparrow MRR	\uparrow R@1	\uparrow R@5	\uparrow R@10	\downarrow MR
Published Results	GNN [12]	52.82	61.37	47.33	77.98	87.83	4.57
	CorefNMN [9]	54.70	61.50	47.55	78.10	88.80	4.40
	RvA [11]	55.59	63.03	49.03	80.40	89.83	4.18
	HACAN [19]	57.17	64.22	50.88	80.63	89.45	4.20
	NMN [9]	58.10	58.80	44.15	76.88	86.88	4.81
	DAN [14]	57.59	63.20	49.63	79.75	89.35	4.30
	DAN † [14]	59.36	64.92	51.28	81.60	90.88	3.92
	ReDAN [15]	61.86	53.13	41.38	66.07	74.50	8.91
	ReDAN+ † [15]	<u>64.47</u>	53.74	42.45	64.68	75.68	6.64
	DualVD [22]	56.32	63.23	49.25	80.23	89.70	4.11
	FGA [13]	56.93	66.22	52.75	82.92	91.08	3.81
	FGA † [13]	57.20	<u>69.30</u>	<u>55.65</u>	<u>86.73</u>	<u>94.05</u>	3.14
	DL-61 [20]	57.32	62.20	47.90	80.43	89.95	4.17
	DL-61 † [20]	57.88	63.42	49.30	80.77	90.68	3.97
	MReal - BDAI * [21]	74.02	52.62	40.03	68.85	79.15	6.76
Leaderboard Entries	LF	45.31	55.42	40.95	72.45	82.83	5.95
	HRE	45.46	54.16	39.93	70.45	81.50	6.41
	MN	47.50	55.49	40.98	72.30	83.30	5.92
	MN-Att	49.58	56.90	42.43	74.00	84.35	5.59
	LF-Att	51.63	60.41	46.18	77.80	87.30	4.75
	MS ConvAI	55.35	63.27	49.53	80.40	89.60	4.15
	USTC-YTH	56.47	61.44	47.65	78.13	87.88	4.65
	UET-VNU	57.40	59.50	45.50	76.33	85.82	5.34
	square	60.16	61.26	47.15	78.73	88.48	4.46
	MS D365 AI	64.47	53.73	42.45	64.68	75.68	6.63
Ours	Random init	60.40	65.53	51.03	83.45	91.83	3.60
	w/ CC [28]+VQA [29]	63.87	67.50	53.85	84.68	93.25	3.32
	CE	74.47	50.74	37.95	64.13	80.00	6.28
	CE + NSP	68.08	63.92	50.78	79.53	89.60	4.28

MR – which depend on the original sparse annotations in VisDial *i.e.* follow-up answers provided in human-human dialog.

We begin by visualizing the distribution of dense relevance scores for these sparse ground-truth (GT) answers in Fig. 2b and observe that $\sim 50\%$ GT answers have relevance ≤ 0.8 , and $\sim 30\%$ have relevance ≤ 0.6 . Thus, there is some degree of misalignment between dense and sparse annotations – answers originally provided during human-human dialog in VisDial were not always judged to be relevant by all humans during the post-hoc dense annotation phase.

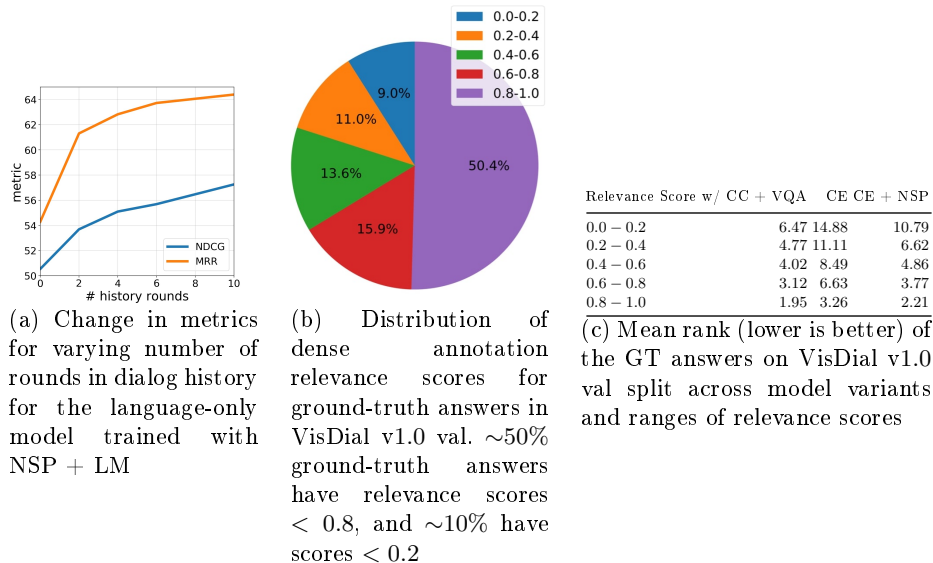


Fig. 2

Why are GT and dense annotations misaligned? We notice that many questions with discrepancy between GT and dense annotations are somewhat subjective. For *e.g.*, in row 1, round 7 (Fig. 5), Q: ‘what color is the chair?’, the GT answer is ‘black’ but the chair is in shadow and it is difficult to accurately identify its color. And thus, we expect to see variance when multiple humans are polled for the answer. Instead, the GT answer is just one sample from the human answer distribution, not necessarily from its peak. In general, the dense annotations seem less wrong than GT (as they are sourced by consensus) since they are safer – often resolving to answers like ‘I cannot tell’ when there is uncertainty / subjectivity – but also uninformative – not conveying additional information *e.g.* ‘I think 3 but they are occluded so it is hard to tell’ – since such nuanced answers are not part of the list of answer options in VisDial [2].

Model performance on GT vs. dense annotations. Table 2c shows mean ranks of these GT answers as predicted by three model variants – ViLBERT w/ CC + VQA, CE, and CE + NSP – grouped by dense relevance scores. The ‘CE’ model gets worse mean ranks than ‘w/ CC + VQA’ for all GT answers, since it is no longer trained with these GT answers during dense annotation finetuning. The CE model assigns low mean ranks to GT answers with higher relevance scores (≥ 0.8), which translates to a high NDCG score (Table 3). But it assigns poor mean ranks to GT answers with relatively lower relevance scores (≤ 0.8), and since ~50% GT answers have relevance scores ≤ 0.8 , this hurts MRR, R@k, MR for the CE model (Table 3).

Next, we consider the top-50 most-relevant answer options (occurring ≥ 10 times) as per dense annotations in VisDial v1.0 val (not restricting ourselves to

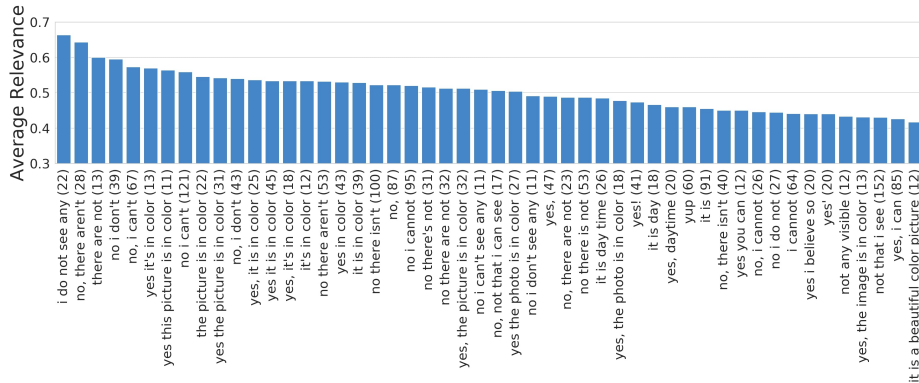


Fig. 3: Mean relevance scores and counts for top-50 most-relevant answers from VisDial v1.0 val dense annotations. These contain several sets of paraphrases – {“yes it’s in color”, “yes this picture is in color”, “the picture is in color”, “yes the picture is in color”, “yes, it is in color”, “yes it is in color”, “yes, it’s in color”, “yes in color”}, *etc.* and have a bias towards binary answers

only GT answers). Fig. 3 shows the mean relevance scores for this set, and Fig. 4 shows the mean ranks assigned to these answers by our models. The CE model gets better mean ranks in this set compared to Base, leading to high NDCG.

Qualitative examples. Finally, we present uniformly sampled example answer predictions on VisDial v1.0 val from our models along with the ground-truth dialog sequences in Fig. 5 and present additional samples in the appendix. In these examples, consistent with the Visual Dialog task definition [2], at every round of dialog, the model gets the image, ground-truth human dialog history (including caption), and follow-up question as input, and predicts the answer. Specifically, the model ranks 100 answer options. Here we show the top-1 prediction.

We make a few observations. 1) The Base model is surprisingly accurate, *e.g.* in row 2, round 1 (Fig. 5), Q: ‘can you see any people?’, predicted answer: ‘part of a person’, in row 2, round 10, Q: ‘anything else interesting about the photo?’, predicted answer: ‘the dog is looking up at the person with his tongue out’. 2) The CE model often answers with generic responses (such as ‘I cannot tell’), especially for questions involving some amount of subjectivity / uncertainty, *e.g.* in row 1, round 7, Q: ‘what color is the chair?’, predicted answer: ‘I cannot tell’ (the chair seems to be in shadow in the image), in row 2, round 7, Q: ‘does the dog look happy?’, predicted answer: ‘I can’t tell’ (subjective question). 3) This also highlights a consequence of misalignment between ground-truth and dense annotations. While the ground-truth answer provides *one* reasonable response for the question asked, it is answerer-specific to quite an extent and there may be other correct answers (annotated in the dense annotations). A negative effect of this misalignment is that when finetuned on dense annotations (CE), the model gets rewarded for generic answers (*e.g.* ‘cannot tell’). While being able to capture and reason about uncertainty is a desirable property models should

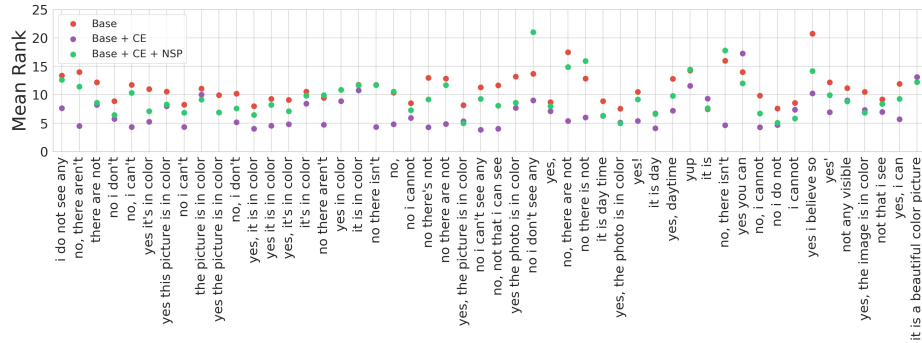


Fig. 4: Predicted mean rank for each of the top-50 most relevant answers as per dense annotations (from Fig. 3) by three model variants – ViLBERT w/ CC + VQA (called ‘Base’), CE, and CE + NSP. The CE model gets lower mean ranks for most answers in this set compared to Base. This leads to significantly higher NDCG, as reported in Table 3 and Table 4, but low MRR, since these relevant answers as per dense annotations do not correlate well with the set of original ground-truth answers, as shown in Fig. 2b

Image + Caption	GT Dialog	Base Model	Base Model + CE	Base Model + CE + NSP
 a chair sitting on the beach with boats in view offshore	Q: are there any people? A: yes Q: is the picture in color? A: yes Q: how many people are there? A: i can't tell Q: what are they doing? A: sitting on the boats Q: is this near a beach? A: yes Q: are they big boats? A: i i Q: what color is the chair? A: white Q: what kind of chair is it? A: beach Q: what is it made out of? A: wood Q: is it a nice day? A: overcast	Q: are there any people? A: no Q: is the picture in color? A: yes Q: how many people are there? A: 1 Q: what are they doing? A: sitting on the boats Q: is this near a beach? A: yes Q: are they big boats? A: no Q: what color is the chair? A: black Q: what kind of chair is it? A: beach Q: what is it made out of? A: metal Q: is it a nice day? A: yes	Q: are there any people? A: no Q: is the picture in color? A: yes, it is Q: how many people are there? A: 1 Q: what are they doing? A: i cannot tell Q: is this near a beach? A: yes Q: are they big boats? A: nope Q: what color is the chair? A: i cannot tell Q: what kind of chair is it? A: i can't tell Q: what is it made out of? A: i can't tell Q: is it a nice day? A: yes	Q: are there any people? A: no Q: is the picture in color? A: yes Q: how many people are there? A: 1 Q: what are they doing? A: sitting on the boats Q: is this near a beach? A: yes Q: are they big boats? A: no Q: what color is the chair? A: black Q: what kind of chair is it? A: beach Q: what is it made out of? A: i can't tell Q: is it a nice day? A: yes
	Q: can you see any people? A: part of a person Q: is the photo in color? A: yes Q: man or woman? A: i can't tell Q: how many dogs? A: just 1 Q: is the dog big or small? A: medium Q: inside or outside? A: hard to tell but i think outside Q: does the dog look happy? A: kind of Q: can you see the ground? A: yes Q: what color is the skateboard? A: black Q: anything else interesting about the photo? A: the dog is looking up at the person with his tongue out	Q: can you see any people? A: part of a person Q: is the photo in color? A: yes Q: man or woman? A: can't see his face Q: how many dogs? A: 1 Q: is the dog big or small? A: big dog Q: inside or outside? A: outside Q: does the dog look happy? A: yes Q: can you see the ground? A: yes Q: what color is the skateboard? A: black Q: anything else interesting about the photo? A: the dog is looking up at the person with his tongue out	Q: can you see any people? A: yes Q: is the photo in color? A: yes Q: man or woman? A: can't tell Q: how many dogs? A: 1 Q: is the dog big or small? A: big dog Q: inside or outside? A: outside Q: does the dog look happy? A: i can't tell Q: can you see the ground? A: yes Q: what color is the skateboard? A: black Q: anything else interesting about the photo? A: no	Q: can you see any people? A: yes Q: is the photo in color? A: yes Q: man or woman? A: can't tell Q: how many dogs? A: 1 Q: is the dog big or small? A: it looks like a large dog Q: inside or outside? A: outside Q: does the dog look happy? A: yes Q: can you see the ground? A: yes Q: what color is the skateboard? A: black Q: anything else interesting about the photo? A: the dog is looking up at the person with his tongue out

Fig. 5: Qualitative samples for three model variants – ViLBERT w/ CC + VQA (called ‘Base’), Base + CE, and Base + CE + NSP

have, it would be more helpful if these agents can convey more information with appropriate qualifiers (*e.g.* ‘I think 3 but they are occluded so it is hard to tell’) than a blanket ‘I cannot tell’. We aim to study this in future work.

7 Implementation

We use the BERT_{BASE} model [35] for the linguistic stream. We use 6 layers of Transformer blocks (with 8 attention heads and a hidden state size of 1024) for the visual stream. The co-attention layers connect the 6 Transformer layers in the visual stream to the last 6 Transformer layers in the linguistic stream. We train on dialog sequences with atmost 256 tokens as most sequences had atmost 256 tokens. During inference, we truncate longer sequences by removing rounds starting from round 1 (we keep the caption). We set all loss coefficients to 1. We use a batch size of 128 for language-only experiments and 80 for other experiments. We use Adam [63] and linearly increase learning rate from 0 to $2e^{-5}$ over $10k$ iterations and decay to $1e^{-5}$ over $200k$ iterations. Our code is available at github.com/vmurahari3/visdial-bert/.

8 Conclusion

We introduce a model for Visual Dialog that enables pretraining on large-scale image-text datasets before transferring and finetuning on VisDial. Our model is an adaptation of ViLBERT [44], and our best single model is pretrained on BooksCorpus [27], English Wikipedia (at the BERT stage), and on Conceptual Captions [28], VQA [29] (at the ViLBERT stage), before finetuning on VisDial, optionally with dense annotations. Our model outperforms prior published results by $> 1\%$ absolute on NDCG and MRR, achieving state-of-the-art results, and providing a simple baseline for future ‘pretrain-then-transfer’ approaches.

Through careful analysis of our results, we find that the recently released dense annotations for the task do not correlate well with the original ground-truth dialog answers, leading to a trade-off when models optimize for metrics that take into account these dense annotations (NDCG) *vs.* the original sparse annotations (MRR). This opens up avenues for future research into better evaluation metrics.

Finally, note that our model is discriminative – it can pick a good answer from a list of answer options – but cannot generate an answer. In the future, we aim to develop robust decoding techniques, based on decoding strategies for transformer-based models introduced in [33, 64], for a strong generative model.

Acknowledgments. The Georgia Tech effort was supported in part by NSF, AFRL, DARPA, ONR YIPs, ARO PECASE, Amazon. AD was supported by fellowships from Facebook, Adobe, Snap Inc. Views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the U.S. Government, or any sponsor.

References

1. H. de Vries, F. Strub, S. Chandar, O. Pietquin, H. Larochelle, and A. Courville, “GuessWhat?! visual object discovery through multi-modal dialogue,” in *CVPR*, 2017. 1, 3, 4
2. A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura, D. Parikh, and D. Batra, “Visual Dialog,” in *CVPR*, 2017. 1, 2, 3, 4, 5, 6, 8, 11, 12
3. F. Strub, H. De Vries, J. Mary, B. Piot, A. Courville, and O. Pietquin, “End-to-end optimization of goal-driven and visually grounded dialogue systems,” *arXiv preprint arXiv:1703.05423*, 2017. 1, 3, 4
4. A. Das, S. Kottur, J. M. Moura, S. Lee, and D. Batra, “Learning Cooperative Visual Dialog Agents with Deep Reinforcement Learning,” in *ICCV*, 2017. 1, 3, 4
5. J. Lu, A. Kannan, J. Yang, D. Parikh, and D. Batra, “Best of both worlds: Transferring knowledge from discriminative learning to a generative visual dialog model,” in *NIPS*, 2017. 1, 3, 4
6. D. Massiceti, N. Siddharth, P. K. Dokania, and P. H. Torr, “Flipdial: A generative model for two-way visual dialogue,” in *CVPR*, 2018. 1, 3, 4
7. Q. Wu, P. Wang, C. Shen, I. Reid, and A. van den Hengel, “Are you talking to me? reasoned visual dialog generation through adversarial learning,” in *CVPR*, 2018. 1, 3, 4
8. U. Jain, S. Lazebnik, and A. G. Schwing, “Two can play this game: visual dialog with discriminative question generation and answering,” in *CVPR*, 2018. 1, 3
9. S. Kottur, J. M. Moura, D. Parikh, D. Batra, and M. Rohrbach, “Visual coreference resolution in visual dialog using neural module networks,” in *ECCV*, 2018. 1, 3, 4, 10
10. S.-W. Lee, T. Gao, S. Yang, J. Yoo, and J.-W. Ha, “Large-scale answerer in questioner’s mind for visual dialog question generation,” in *ICLR*, 2019. 1, 3, 4
11. Y. Niu, H. Zhang, M. Zhang, J. Zhang, Z. Lu, and J.-R. Wen, “Recursive visual attention in visual dialog,” in *CVPR*, 2019. 1, 3, 4, 10
12. Z. Zheng, W. Wang, S. Qi, and S.-C. Zhu, “Reasoning visual dialogs with structural and partial observations,” in *CVPR*, 2019. 1, 3, 10
13. I. Schwartz, S. Yu, T. Hazan, and A. G. Schwing, “Factor graph attention,” in *CVPR*, 2019. 1, 3, 10
14. G.-C. Kang, J. Lim, and B.-T. Zhang, “Dual attention networks for visual reference resolution in visual dialog,” in *EMNLP*, 2019. 1, 3, 10
15. Z. Gan, Y. Cheng, A. E. Kholy, L. Li, J. Liu, and J. Gao, “Multi-step reasoning via recurrent dual attention for visual dialog,” in *ACL*, 2019. 1, 3, 10
16. S. Kottur, J. M. Moura, D. Parikh, D. Batra, and M. Rohrbach, “Clevr-dialog: A diagnostic dataset for multi-round reasoning in visual dialog,” in *NAACL*, 2019. 1, 3
17. V. Murahari, P. Chattopadhyay, D. Batra, D. Parikh, and A. Das, “Improving generative visual dialog by answering diverse questions,” in *EMNLP*, 2019. 1, 3, 4
18. R. Shekhar, A. Venkatesh, T. Baumgärtner, E. Bruni, B. Plank, R. Bernardi, and R. Fernández, “Beyond task success: A closer look at jointly learning to see, ask, and guesswhat,” in *NAACL*, 2019. 1, 3
19. T. Yang, Z.-J. Zha, and H. Zhang, “Making history matter: Gold-critic sequence training for visual dialog,” *arXiv preprint arXiv:1902.09326*, 2019. 1, 3, 10
20. D. Guo, C. Xu, and D. Tao, “Image-question-answer synergistic network for visual dialog,” in *CVPR*, 2019. 1, 3, 10

21. J. Qi, Y. Niu, J. Huang, and H. Zhang, “Two causal principles for improving visual dialog,” *arXiv preprint arXiv:1911.10496*, 2019. 1, 3, 10
22. X. Jiang, J. Yu, Z. Qin, Y. Zhuang, X. Zhang, Y. Hu, and Q. Wu, “DualVD: An adaptive dual encoding model for deep visual understanding in visual dialogue,” in *AAAI*, 2020. 1, 3, 10
23. H. Alamri, V. Cartillier, A. Das, J. Wang, A. Cherian, I. Essa, D. Batra, T. K. Marks, C. Hori, P. Anderson, S. Lee, and D. Parikh, “Audio visual scene-aware dialog,” in *CVPR*, 2019. 1
24. H. de Vries, K. Shuster, D. Batra, D. Parikh, J. Weston, and D. Kiela, “Talk the walk: Navigating new york city through grounded dialogue,” *arXiv preprint arXiv:1807.03367*, 2018. 1
25. K. Nguyen and H. Daumé III, “Help, anna! visual navigation with natural multimodal assistance via retrospective curiosity-encouraging imitation learning,” in *EMNLP*, 2019. 1
26. J. Thomason, M. Murray, M. Cakmak, and L. Zettlemoyer, “Vision-and-dialog navigation,” 2019. 1
27. Y. Zhu, R. Kiros, R. Zemel, R. Salakhutdinov, R. Urtasun, A. Torralba, and S. Fidler, “Aligning books and movies: Towards story-like visual explanations by watching movies and reading books,” in *ICCV*, 2015. 2, 7, 8, 14
28. P. Sharma, N. Ding, S. Goodman, and R. Soricut, “Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning,” in *ACL*, 2018. 2, 3, 4, 5, 7, 9, 10, 14
29. S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, “VQA: Visual Question Answering,” in *ICCV*, 2015. 2, 3, 4, 5, 7, 9, 10, 14
30. K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *CVPR*, 2016. 2, 3, 5
31. K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *ICLR*, 2015. 2, 3
32. A. Krizhevsky, I. Sutskever, and G. Hinton, “ImageNet Classification with Deep Convolutional Neural Networks,” in *NIPS*, 2012. 2, 3
33. A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, “Improving language understanding with unsupervised learning,” 2018. 2, 3, 4, 14
34. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *NIPS*, 2017. 2, 3, 4, 7
35. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *NAACL*, 2019. 2, 3, 4, 5, 7, 14
36. Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “RoBERTa: A robustly optimized BERT pretraining approach,” *arXiv preprint arXiv:1907.11692*, 2019. 2, 3, 4
37. Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, “ALBERT: A lite BERT for self-supervised learning of language representations,” *arXiv preprint arXiv:1909.11942*, 2019. 2, 3, 4
38. Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. Salakhutdinov, and Q. V. Le, “Xlnet: Generalized autoregressive pretraining for language understanding,” *arXiv preprint arXiv:1906.08237*, 2019. 2, 3, 4
39. C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *arXiv preprint arXiv:1910.10683*, 2019. 2, 3, 4

40. Y. Zhang, S. Sun, M. Galley, Y.-C. Chen, C. Brockett, X. Gao, J. Gao, J. Liu, and B. Dolan, “DialogPT: Large-scale generative pre-training for conversational response generation,” *arXiv preprint arXiv:1911.00536*, 2019. 2, 3, 4
41. O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, *et al.*, “Imagenet large scale visual recognition challenge,” *IJCV*, 2015. 2, 4
42. R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, *et al.*, “Visual genome: Connecting language and vision using crowdsourced dense image annotations,” *IJCV*, 2017. 2, 5
43. A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, “GLUE: A multi-task benchmark and analysis platform for natural language understanding,” *arXiv preprint arXiv:1804.07461*, 2018. 2, 4
44. J. Lu, D. Batra, D. Parikh, and S. Lee, “ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks,” in *NeurIPS*, 2019. 2, 3, 4, 5, 14
45. L. H. Li, M. Yatskar, D. Yin, C.-J. Hsieh, and K.-W. Chang, “VisualBERT: A simple and performant baseline for vision and language,” *arXiv preprint arXiv:1908.03557*, 2019. 2, 4
46. H. Tan and M. Bansal, “LXMERT: Learning cross-modality encoder representations from transformers,” *arXiv preprint arXiv:1908.07490*, 2019. 2, 4
47. Y.-C. Chen, L. Li, L. Yu, A. E. Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, “UNITER: Learning UNiversal Image-TEText Representations,” *arXiv preprint arXiv:1909.11740*, 2019. 2, 4
48. G. Li, N. Duan, Y. Fang, D. Jiang, and M. Zhou, “Unicoder-VL: A universal encoder for vision and language by cross-modal pre-training,” *arXiv preprint arXiv:1908.06066*, 2019. 2, 4
49. W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, “VL-BERT: Pre-training of generic visual-linguistic representations,” *arXiv preprint arXiv:1908.08530*, 2019. 2, 4
50. C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid, “VideoBERT: A joint model for video and language representation learning,” *arXiv preprint arXiv:1904.01766*, 2019. 2, 4
51. R. Zellers, Y. Bisk, A. Farhadi, and Y. Choi, “From recognition to cognition: Visual commonsense reasoning,” in *CVPR*, 2019. 2, 4
52. A. Suhr, S. Zhou, A. Zhang, I. Zhang, H. Bai, and Y. Artzi, “A corpus for reasoning about natural language grounded in photographs,” in *ACL*, 2019. 2, 4
53. N. Xie, F. Lai, D. Doran, and A. Kadav, “Visual entailment: A novel task for fine-grained image understanding,” *arXiv preprint arXiv:1901.06706*, 2019. 2, 4
54. P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, “From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions,” *TACL*, 2014. 2, 4
55. K.-H. Lee, X. Chen, G. Hua, H. Hu, and X. He, “Stacked cross attention for image-text matching,” in *ECCV*, 2018. 2, 4
56. S. Kazemzadeh, V. Ordonez, M. Matten, and T. L. Berg, “ReferItGame: Referring to Objects in Photographs of Natural Scenes,” in *EMNLP*, 2014. 2, 4
57. W. Hao, C. Li, X. Li, L. Carin, and J. Gao, “Towards learning a generic agent for vision-and-language navigation via pre-training,” in *CVPR*, 2020. 2, 4
58. A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch:

- An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32* (H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, eds.), pp. 8024–8035, Curran Associates, Inc., 2019. [3](#)
59. T. Wolf, V. Sanh, J. Chaumond, and C. Delangue, “TransferTransfo: A transfer learning approach for neural network based conversational agents,” *arXiv preprint arXiv:1901.08149*, 2019. [5](#)
 60. P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, “Bottom-up and top-down attention for image captioning and visual question answering,” *arXiv preprint arXiv:1707.07998*, 2017. [5](#)
 61. S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” in *NIPS*, 2015. [5](#)
 62. T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft COCO: Common Objects in Context,” in *ECCV*, 2014. [6](#)
 63. D. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization,” in *ICLR*, 2015. [14](#)
 64. M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, “Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension,” *arXiv preprint arXiv:1910.13461*, 2019. [14](#)