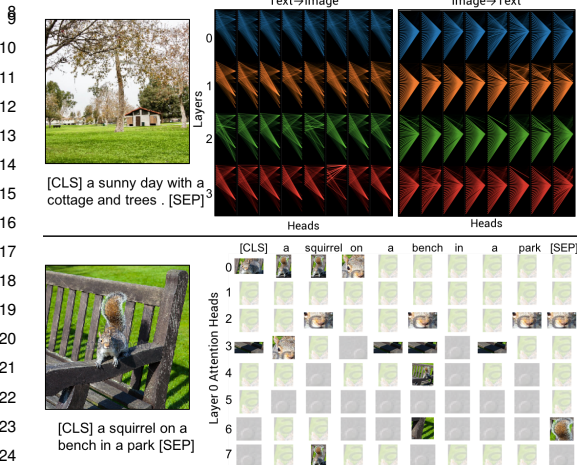
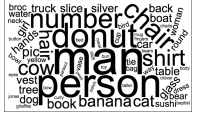


1 We thank the reviewers for the thoughtful feedback! We are encouraged that all voted to accept, finding the paper
 2 clear / well-organized [R1]; our approach “very interesting” [R3] and novel [R2 R3]; our results significant and
 3 well-demonstrated [R1 R2]; and likely to be built on by the community [R1]. We are pleased they recognized the value
 4 of transferring visio-linguistic pretraining [R1 R2 R3] and the demonstrated benefits of our co-attentional two-stream
 5 model over a direct extension of BERT [R2 R3]. We respond to select comments below but will address all feedback.
 6 **Improved performance.** After submission, refined LR schedules raised performance across all tasks – passing the
 7 recent VQA challenge winner, setting a new SoTA. Will update paper with details and release the codebase if accepted.

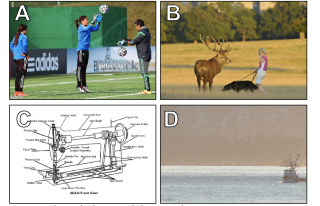


[R1] Visualization of coattention over multimodal inputs. Visualizing BERT-style models is an open research area; we applied the method of [Fig. arXiv 2019] and observed some trends – providing representative examples on the left. **[Top]** We show attention for each layer (rows) and head (cols) with attention focus (colored lines) shown going from source to target (left-to-right), text→image co-attention tends to be better grounded in early layers before converging to a set of somewhat arbitrary regions as depth increases – see attention focus concentrate more over layers. In contrast, image→text co-attention often focuses on the high-level sentence representation in the SEP token already developed on the text side early on, but spreads out somewhat later. **[Bottom]** We also show the most attended patch for each attention head for each word in the first layer for an example. Many heads focus on a small set of “default” patches (faded for clarity); but, the noun phrases surrounding “squirrel” and “bench” focus more on relevant regions.

25 **[R1] Additional result analysis.** We investigate the RefCOCO+ task. For each noun occurring
 26 in a referring expression, we counted the number of instances where ViLBERT (full) succeeded
 27 and ViLBERT (w/o pretrain) fails (and vice versa). The wordcloud on the right shows those nouns
 28 with the highest performance delta. We will perform more task specific task in supplementary.



29 **[R2] All tasks are “image captioning (or closely related)” so this is “effectively**
 30 **transfer learning from a large captioning dataset to a small one.”** We respectfully
 31 disagree. Due to its automatic collection from the web, Conceptual Captions (CC) is
 32 fairly distinct from curated vision-and-language datasets (examples right). Even for
 33 the closely-related caption-based image retrieval task, it was not obvious to us that this
 34 weakly-aligned web data would help. Further, our other transfer tasks differ significantly
 35 from CC. VQA and VCR both ask grounded questions like “Is there something to cut
 36 the vegetables with?” (VQA) This is not caption-like and requires reasoning (knives
 37 cut) and grounding. VCR extends to answer justifications like “[Person3] is delivering
 38 food to the table, and she might not know whose order is whose” that often refer to
 39 actions and intentions of individuals. Referring expressions focus on aligning small image
 40 text like “guy in yellow dribbling ball” – both being quite different from whole-image
 41 descriptive captioning. *However, a common need for visual grounding underpins these tasks and is precisely what we target with our pretraining strategy.*



A. exercises during a training session
 B. deer white walking a dog
 C. diagram of a modern sowing machine
 D. fishing boat returning to port in winter in mid afternoon, a frigid breeze giving lie to the warm glowing light

42 **[R2] Additional pretraining ablations.** Great suggestions! We report separate image-text
 43 pretraining (w/o corr – masking loss only and zeroed co-attn), without alignment loss (w/o
 44 align), and without masking loss (w/o mask) ablations to the right (only two tasks due to
 45 time). All ablations degrade performance – especially w/o mask which struggles to train
 46 downstream tasks. These ablations are valuable and will be added to the paper.

Method	VQA	RefCOCO+
full	66.59	70.38
w/o corr	64.85	68.04
w/o align	64.61	68.49
w/o mask	42.43	10.00

47 **[R3] Given the use of Conceptual Captions (CC), are the comparisons to baselines**
 48 **fair?** We believe these comparisons are fair. We agree that CC is a large, additional data source; however, being able to
 49 leverage this additional data for a diverse range of vision and language tasks is precisely our contribution! Existing
 50 approaches to vision and language tasks are simply not designed to do so – for instance, it is unclear how to train a
 51 standard VQA model like BAN with CC captioning data. Arguing from analogy, the widespread transfer of deep models
 52 pretrained on ImageNet also leveraged more data during pretraining; however, we do not find it unfair to pre-deep
 53 learning approaches that were not equipped to leverage that data. Finally, note that unlike ImageNet, CC is webly
 54 supervised, and did not involve expensive human annotation. **We acknowledge that in caption-based image retrieval,**
 55 **CC data could have been used to pretrain existing work for a more direct architectural comparison – we will address.**

56 **[R3] “If I understood correctly, [the w/o pretrain model] does not use any visual features.** That is not the case.
 57 Like all our models, the “w/o pretrain” model is initialized from a trained visual feature extractor (Faster RCNN) and
 58 language model (BERT). We use “w/o pretrain” to note that the model has not undergone our visio-linguistic pretraining
 59 on the CC dataset (L264). To reduce confusion, we will use “w/o grounding pretraining” and clarify relevant sentences.