# Supplementary material: A Sentence is Worth a Thousand Pixels

Sanja Fidler
TTI Chicago
fidler@ttic.edu

Abhishek Sharma
University of Maryland
bhokaal@cs.umd.edu

Raquel Urtasun
TTI Chicago
rurtasun@ttic.edu

In the supplementary material for paper [3] we provide more quantitative and qualitative results. Table 1 shows the results obtained with our approach when using GT cardinality potentials instead of text cardinality potentials. Note that GT cardinality only affects the potential on $z$ and potential on the number of detection boxes. The remaining potentials (along with the score of the detector) are kept the same. "GT noneg" denotes the experiment, where we encourage at least as many boxes as dictated by the cardinality to be on in the image. With "GT - neg" we denote the experiment where we also **suppress** the boxes for classes with cardinality $0$. This means that for images where GT cardinality for a class is $0$, we simply do not input any boxes of that class into the model.

Interestingly, by using GT instead of extracted cardinalities from text, we do not observe a significant boost in performance. This means that our holistic model is able to (close to) fully exploit the available information about the scene.

Note, that conditioned on whether a noun is mentioned in text for a particular image, our approach uses additional boxes which originally did not pass the detector's threshold. In Table 2, we compare our approach with that of Yao et al. [5], where we use as many boxes as used in our approach. To do this, we reduce the thresholds of the detector, so that for each class the number of boxes exactly matches the number of boxes that our approach uses for the class. Notice that the original approach improved its performance by only $0.4\%$, while our performance is much higher ($12.5\%$ improvement). This means that the success of our approach is not due to an increased number of used boxes, but in how they are used in the model by exploiting additional text information.

In the future we plan to incorporate more powerful segmentation features [1] and detector [2].

| | back. | aerop. | bicycle | bird | boat | bottle | bus | car | cat | chair | cow | dtable | dog | horse | mbike | person | pplant | sheep | sofa | train | monitor | **averg.** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Oracle Z - noneg | 76.8 | 36.7 | 28.3 | 34.8 | 21.9 | 30.9 | 56.1 | 47.6 | 36.8 | 10.0 | 58.2 | 28.8 | 33.4 | 54.8 | 42.6 | 41.8 | 15.1 | 28.1 | 16.3 | 35.7 | 48.7 | 37.3 |
| Oracle Z - neg | 76.8 | 31.2 | 28.2 | 34.7 | 21.7 | 31.1 | 56.0 | 50.3 | 36.7 | 10.5 | 57.4 | 29.5 | 34.4 | 55.1 | 42.5 | 41.9 | 16.9 | 32.2 | 18.8 | 35.7 | 49.2 | 37.7 |
| ours | 76.9 | 31.3 | 29.7 | 37.3 | 27.7 | 29.5 | 52.1 | 40.0 | 38.0 | 6.6 | 55.9 | 25.2 | 33.2 | 38.2 | 44.3 | 42.5 | 15.2 | 32.0 | 20.2 | 40.7 | 48.4 | 36.4 |

Table 1. Comparison to oracle Z (see text for details).

| | back. | aerop. | bicycle | bird | boat | bottle | bus | car | cat | chair | cow | dtable | dog | horse | mbike | person | pplant | sheep | sofa | train | monitor | **averg.** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Textonboost (unary) [4] | 77.8 | 14.1 | 3.4 | 0.7 | 11.3 | 3.3 | 25.5 | 30.9 | 10.3 | 0.7 | 13.2 | 10.8 | 5.2 | 15.1 | 31.8 | 41.0 | 0.0 | 3.7 | 2.4 | 17.1 | 33.7 | 16.8 |
| Holistic Scene Understanding [5] | 77.3 | 25.6 | 12.9 | 14.2 | 19.2 | 31.0 | 34.6 | 38.6 | 16.1 | 7.4 | 11.9 | 9.0 | 13.9 | 25.4 | 31.7 | 38.1 | 11.2 | 18.8 | 6.2 | 23.6 | 34.4 | 23.9 |
| [5] num boxes from text | 77.8 | 26.7 | 14.3 | 11.5 | 18.6 | 30.8 | 34.4 | 37.9 | 17.2 | 5.7 | 19.0 | 7.3 | 12.4 | 27.3 | 36.5 | 37.1 | 11.6 | 9.4 | 6.2 | 25.7 | 43.8 | 24.3 |
| ours | 76.9 | 31.3 | 29.7 | 37.3 | 27.7 | 29.5 | 52.1 | 40.0 | 38.0 | 6.6 | 55.9 | 25.2 | 33.2 | 38.2 | 44.3 | 42.5 | 15.2 | 32.0 | 20.2 | 40.7 | 48.4 | 36.4 |

Table 2. Comparison to the state-of-the-art that utilizes only image information in the UIUC sentence dataset. By leveraging text information our approach improves 12.5% AP. Note that this dataset contains only 600 PASCAL VOC 2008 images for training, and thus is significantly a more difficult task than recent VOC challenges which have up to 10K training images.

# References

[1] J. Carreira, R. Caseiroa, J. Batista, and C. Sminchisescu. Semantic segmentation with second-order pooling. In *ECCV*, 2012. 1

[2] S. Fidler, R. Mottaghi, A. Yuille, and R. Urtasun. Bottom-up segmentation for top-down detection. In *CVPR*, 2013. 1

[3] S. Fidler, A. Sharma, and R. Urtasun. A sentence is worth a thousand pixels. In *CVPR*, 2013. 1

[4] J. Shotton, M. Johnson, and R. Cipolla. Semantic texton forests for image categorization and segmentation. In *CVPR*, 2008. 1

[5] Y. Yao, S. Fidler, and R. Urtasun. Describing the scene as a whole: Joint object detection, scene classification and semantic segmentation. In *CVPR*, 2012. 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12

In Figure 1- 8 we show some examples of successful results, while Figure 9- 11 shows (partial) failures.
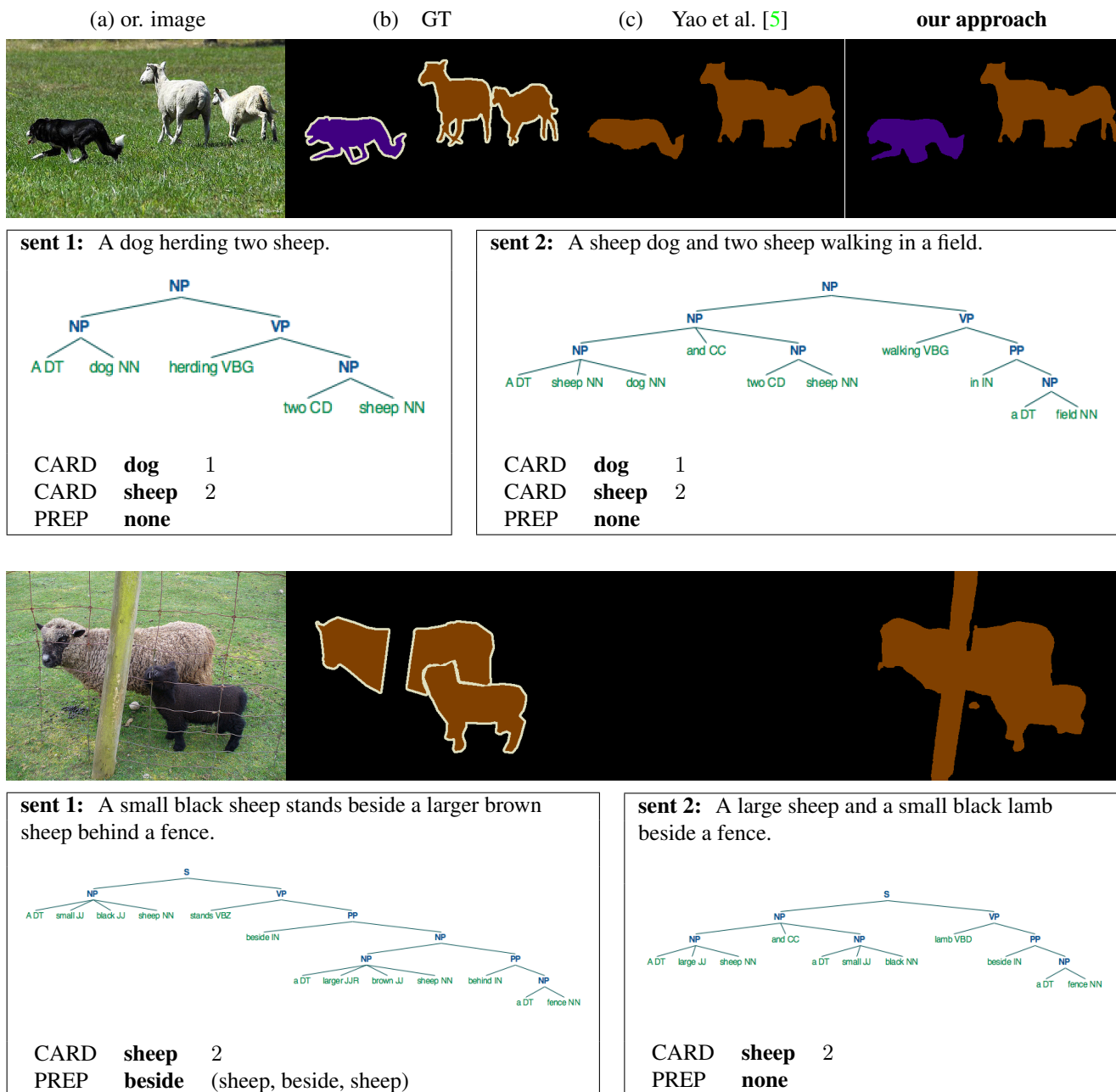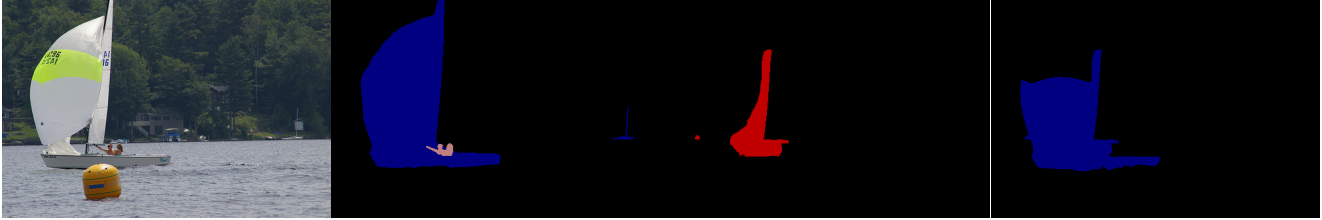
| (a) or. image | (b)  GT | (c)  Yao et al. [5] | **our approach** |



**sent 1:** A dog herding two sheep.



CARD    **dog**     1
CARD    **sheep**   2
PREP    **none**

**sent 2:** A sheep dog and two sheep walking in a field.



CARD    **dog**     1
CARD    **sheep**   2
PREP    **none**



**sent 1:** A small black sheep stands beside a larger brown sheep behind a fence.



CARD    **sheep**   2
PREP    **beside**  (sheep, beside, sheep)

**sent 2:** A large sheep and a small black lamb beside a fence.



CARD    **sheep**   2
PREP    **none**

Figure 1. Examples of results. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.

| (a) or. image | (b) GT | (c) Yao et al. [5] | **our approach** |

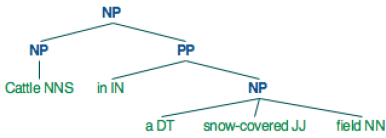**sent 1:** Sail boat on water with two people riding inside.

CARD **boat** 1
CARD **person** 2
PREP **none**

**sent 2:** Small sailboat with spinnaker passing a buoy.

CARD **boat** 1
PREP **none**

**sent 1:** A middle eastern couple sitting on a couch holding their baby and displaying a gift.

CARD **person** 2
CARD **sofa** 1
PREP **on** (person, on, sofa)

**sent 2:** Couple with newborn baby.

CARD **person** 2
PREP **near** (person, with, person)

Figure 2. Examples of results. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.

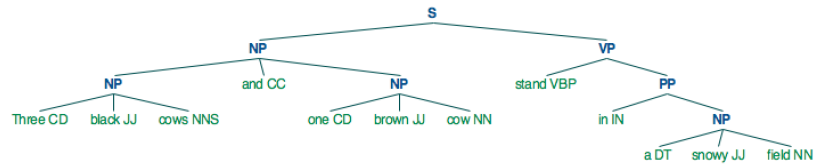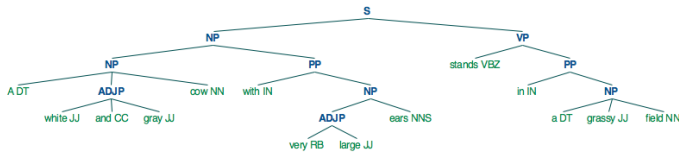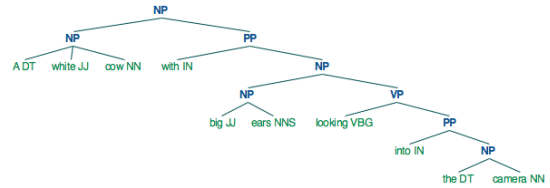(a) or. image     (b)   GT     (c)   Yao et al. [5]     **our approach**

**sent 1:** Passengers loading onto a train with a green and black steam engine.

S — VP — Passengers VBZ — VP — loading VBG — PP — onto IN — NP — a DT — train NN — PP — with IN — NP — a DT — ADJP — green JJ — and CC — black JJ — steam NN

| CARD | **person** | 2 |
| CARD | **train** | 2 |
| PREP | **none** | |

**sent 2:** Several people waiting to board the train.

S — NP — NP — Several JJ — people NNS — VP — waiting VBG — PP — to TO — NP — board NN — the DT — VP — train VBP

| CARD | **person** | 2 |
| CARD | **train** | 1 |
| PREP | **none** | |

**sent 1:** A table is set with wine and dishes for two people.

S — NP — A DT — table NN — VP — is VBZ — set VBN — VP — PP — with IN — NP — wine NN — and CC — dishes NNS — PP — for IN — NP — two CD — people NNS

| CARD | **diningtable** | 1 |
| CARD | **person** | 2 |
| PREP | **none** | |

**sent 2:** A table set for two.

S — NP — A DT — table NN — VP — set VBD — PP — for IN — NP — two CD

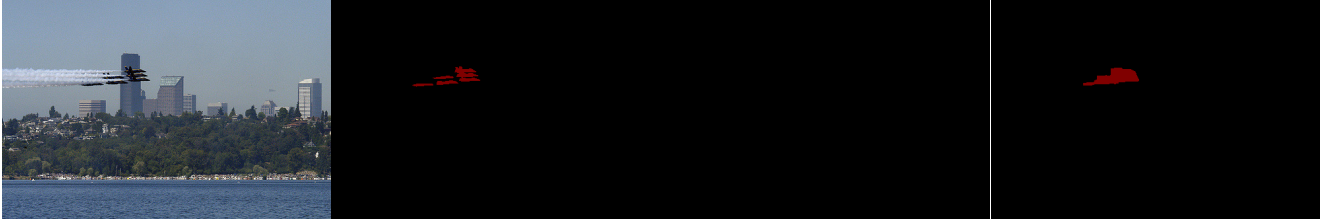| CARD | **diningtable** | 1 |
| PREP | **none** | |

Figure 3. Examples of results. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.
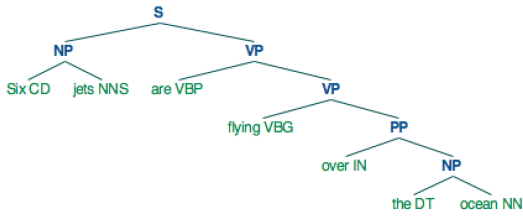
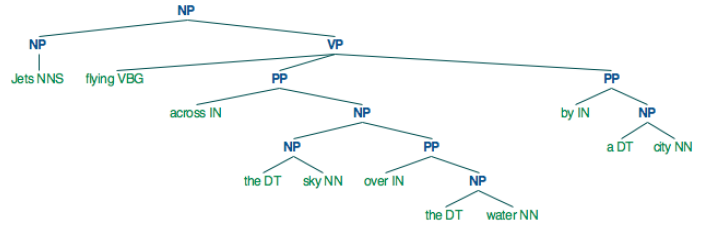(a) or. image     (b) GT     (c) Yao et al. [5]     **our approach**

**sent 1:** Cattle in a snow-covered field.

NP
NP — PP
Cattle NNS   in IN   NP
a DT   snow-covered JJ   field NN

| CARD | **none** | |
| PREP | **none** | |

**sent 2:** Three black cows and one brown cow stand in a snowy field.

S
NP — VP
NP — and CC — NP — stand VBP — PP
Three CD   black JJ   cows NNS    one CD   brown JJ   cow NN    in IN   NP
a DT   snowy JJ   field NN

| CARD | **cow** | 5 |
| PREP | **none** | |

**sent 1:** An old fashioned passenger bus with open windows.

S
NP — VP
An DT   old JJ   fashioned VBN   NP — PP
passenger NN   bus NN   with IN   NP
open JJ   windows NNS

| CARD | **bus** | 1 |
| PREP | **none** | |

**sent 2:** The idle tourist bus awaits its passengers.

S
NP — VP
The DT   idle JJ   tourist NN   bus NN   awaits VBZ   NP
its {VBZ aw}   passengers NNS

| CARD | **bus** | 1 |
| CARD | **person** | 2 |
| PREP | **none** | |

Figure 4. Examples of results. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.
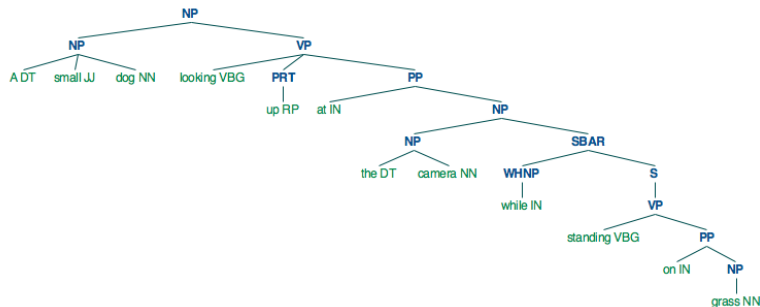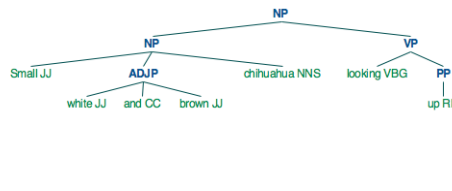
| (a) or. image | (b) GT | (c) Yao et al. [5] | **our approach** |

**sent 1:** A white and gray cow with very large ears stands in a grassy field.

| CARD | **cow** | 1 |
| PREP | **none** | |

**sent 2:** A white cow with big ears looking into the camera.

| CARD | **cow** | 1 |
| PREP | **none** | |

**sent 1:** An orange tabby cat sleeping on the sofa.

| CARD | **cat** | 1 |
| CARD | **sofa** | 1 |
| PREP | **on** | (cat, on, sofa) |

**sent 2:** Yellow striped cat resting on blue sofa.

| CARD | **cat** | 1 |
| CARD | **sofa** | 1 |
| PREP | **on** | (cat, on, sofa) |

Figure 5. Examples of results. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.
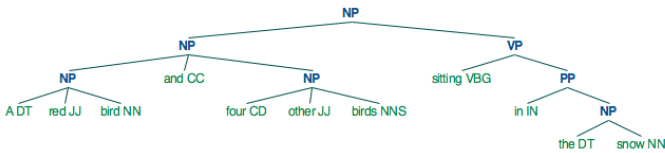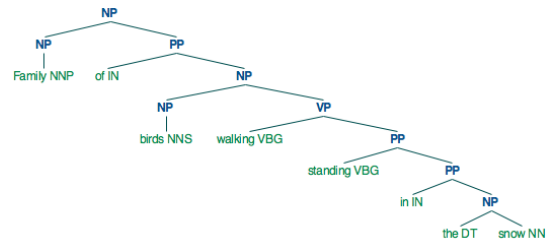
6

| (a) or. image | (b)  GT | (c)  Yao et al. [5] | **our approach** |



**sent 1:**  Six jets are flying over the ocean.

S
→ NP
→→ Six CD   jets NNS
→ VP
→→ are VBP
→→ VP
→→→ flying VBG
→→→ PP
→→→→ over IN
→→→→ NP
→→→→→ the DT   ocean NN

CARD     **aeroplane**     6
PREP     **none**

**sent 2:**  Jets flying across the sky over the water by a city.

NP
→ NP
→→ Jets NNS   flying VBG
→ VP
→→ PP
→→→ across IN
→→→ NP
→→→→ NP
→→→→→ the DT   sky NN
→→→→ PP
→→→→→ over IN
→→→→→ NP
→→→→→→ the DT   water NN
→→ PP
→→→ by IN
→→→ NP
→→→→ a DT   city NN

CARD     **aeroplane**     2
PREP     **none**



**sent 1:**  Dark colored dog standing in grassy field.

NP
→ NP
→→ Dark JJ   colored VBN   dog NN   standing NN   in IN
→ PP
→→ in IN
→→ NP
→→→ grassy JJ   field NN

CARD     **dog**     1
PREP     **none**

**sent 2:**  A grey dog standing in a grassy field.

NP
→ NP
→→ A DT   gray JJ   dog NN   standing NN   in IN
→ PP
→→ in IN
→→ NP
→→→ a DT   grassy JJ   field NN

CARD     **dog**     1
PREP     **none**

Figure 6. Examples of results. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.
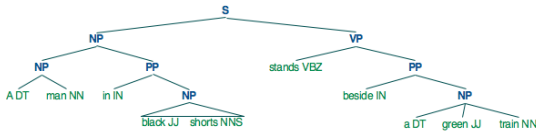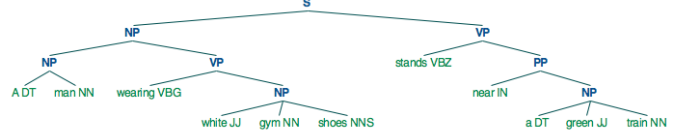
| (a) or. image | (b) GT | (c) Yao et al. [5] | **our approach** |

**sent 1:** A small dog looking up at the camera while standing on grass.

CARD **dog** 1
PREP **none**

**sent 2:** Small white and brown chihuahua looking up.

CARD **dog** 1
PREP **none**

**sent 1:** This is an image of the back of some old televisions.

CARD **tvmonitor** 2
PREP **none**

**sent 2:** Three televisions, on on the floor, the other two on a box.

CARD **tvmonitor** 3
PREP **none**

Figure 7. Examples of results. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.

8

| (a) or. image | (b) GT | (c) Yao et al. [5] | **our approach** |

**sent 1:** A red bird and four other birds sitting in the snow.

CARD  **bird**  5
PREP  **none**

**sent 2:** Family of birds walking standing in the snow.

CARD  **bird**  2
PREP  **none**

**sent 1:** A large bird is flying through the air.

CARD  **bird**  1
PREP  **none**

**sent 2:** A large bird of prey flying towards the camera.

CARD  **bird**  1
PREP  **none**

Figure 8. Examples of results. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.
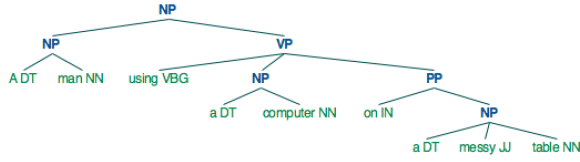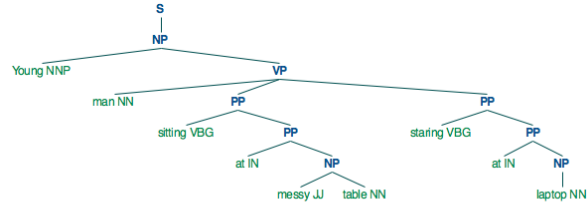
| (a) or. image | (b) GT | (c) Yao et al. [5] | **our approach** |

**sent 1:** A man in black shorts stands beside a green train.



| CARD | **person** | 1 |
| CARD | **train** | 1 |
| PREP | **near** | (person, beside, train) |

**sent 2:** A man wearing white gym shoes stands near a green train.



| CARD | **person** | 1 |
| CARD | **train** | 1 |
| PREP | **near** | (person, near, train) |

**sent 1:** Older woman in kitchen with her two domestic cats.



| CARD | **cat** | 1 |
| CARD | **person** | 1 |
| PREP | **near** | (person, with, cat) |

**sent 2:** The old lady is standing in the kitchen with two cats at her feet.



| CARD | **cat** | 1 |
| CARD | **person** | 1 |
| PREP | **near** | (person, with, cat) |

Figure 9. Examples of (partial) failures. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.

(a) or. image      (b) GT      (c) Yao et al. [5]      **our approach**

**sent 1:** A man using a computer on a messy table.

| CARD | diningtable | 1 |
| CARD | person | 1 |
| CARD | tvmonitor | 1 |
| PREP | on | (person, on, diningtable) |

**sent 2:** Young man sitting at messy table staring at laptop.

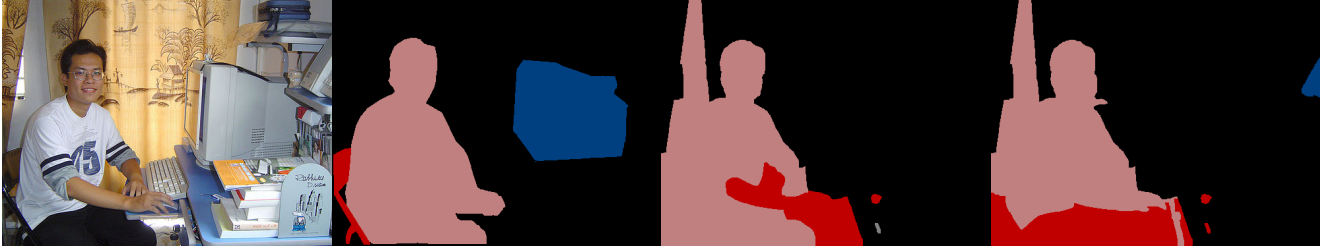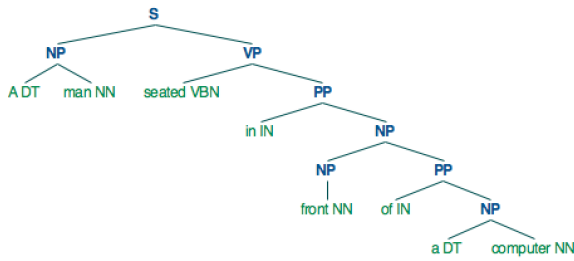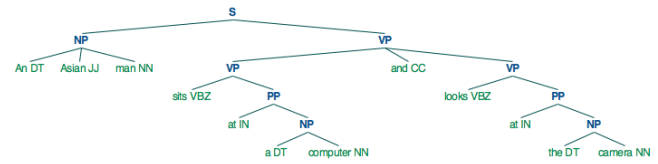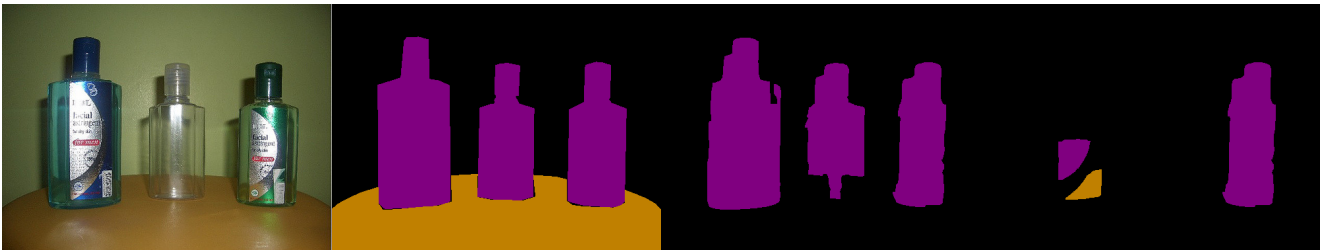| CARD | diningtable | 1 |
| CARD | person | 1 |
| CARD | tvmonitor | 1 |
| PREP | none | |

**sent 1:** Older woman dressed in white night gown on tan couch.

| CARD | person | 1 |
| CARD | sofa | 1 |
| PREP | on | (person, on, sofa) |

**sent 2:** Woman in a floral nightgown sits on a couch.

| CARD | person | 1 |
| CARD | sofa | 1 |
| PREP | on | (person, on, sofa) |

Figure 10. Examples of (partial) failures. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.

Figure 11. Examples of (partial) failures. We show the original image, GT segmentation, the output of [5], and our result. Below the images we show the exemplar sentences with parse tree, with extracted object classes, cardinalities and prepositions.