

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

## Appendix for “End-to-End Instance Segmentation with Recurrent Attention”

Anonymous CVPR submission

Paper ID 3051

### A. Training procedure specification

We used the Adam optimizer [2] with learning rate 0.001 and batch size of 8. The learning rate is multiplied by 0.85 for every 5000 steps of training.

#### A.1. Scheduled sampling

We denote  $\theta_t$  as the probability of feeding in ground-truth segmentation that has the greatest overlap with the previous prediction, as opposed to model output.  $\theta_t$  decays exponentially as training proceeds, and for larger  $t$ , the decay occurs later:

$$\theta_t = \min \left( \Gamma_t \exp \left( -\frac{\text{epoch} - S}{S_2} \right), 1 \right) \quad (1)$$

$$\Gamma_t = 1 + \log(1 + Kt) \quad (2)$$

where  $\text{epoch}$  is the training index,  $S$ ,  $S_2$ , and  $K$  are constants. In the experiments reported here, these values are 10000, 2885, and 3.

### B. Evaluation metrics

We include the details of evaluation metrics here. Symmetric best dice (SBD) (Eq. 3-5) is used on the CVPPP dataset. Mean (un)weighted coverage (MUCov, MW Cov) (Eq. 6-7) is used on the KITTI dataset. Average precision (AP) (Eq. 9) is used on the Cityscapes dataset.

$$\text{DICE}(A, B) = \frac{2|A \hat{B}|}{|A| + |B|} \quad (3)$$

$$\text{BD}(\{A_i\}, B) = \max_i \text{DICE}(A_i, B) \quad (4)$$

$$\begin{aligned} \text{SBD}(y_i, \{y_j^*\}) &= \min \left( \frac{1}{N} \sum_j \text{BD}(\{y_i\}, y_j^*), \right. \\ &\quad \left. \frac{1}{M} \sum_i \text{BD}(\{y_j^*\}, y_i) \right) \end{aligned} \quad (5)$$

**Table 1:** MS-COCO Zebra Results

	MWCov $\uparrow$	MUCov $\uparrow$	$ \text{DiC}  \downarrow$	Acc. $\uparrow$
detect [1]	-	-	2.56	-
aso-sub [1]	-	-	1.03	-
Ours	69.2	64.2	<b>0.79</b>	0.57

$$\text{MUCov}(\{y_i\}, \{y_j^*\}) = \sum_i \frac{1}{N} \max_j \text{IoU}(y_i, y_j^*) \quad (6)$$

$$\text{MWCov}(\{y_i\}, \{y_j^*\}) = \sum_i w_{\text{cov},i} \max_j \text{IoU}(y_i, y_j^*) \quad (7)$$

$$w_{\text{cov},i} = \frac{|y_i|}{\sum_i |y_i|} \quad (8)$$

$$\begin{aligned} \text{AP}(\{y_i\}, \{y_j^*\}) &= \max_s \sum_\theta \sum_j Pr(y_{s(i)}, y_j) \cdot \\ &\quad \mathbb{1}[\text{IoU}(y_{s(i)}, y_j^*) \geq \theta], \end{aligned} \quad (9)$$

### C. More experimental results

We include the segmentation and counting performance on the MS-COCO zebra images in Table 1. In terms of counting, our model out-performs a baseline method that runs an object detector and then non-maximal suppression, and a new associative-subitizing method [1].

### D. Model architecture

#### D.1. Foreground + Orientation FCN

We resize the image to uniform size. For CVPPP and MS-COCO dataset, we adopt a uniform size of  $224 \times 224$ , for KITTI, we adopt  $128 \times 448$ , and for Cityscapes  $256 \times 512$  (4x downsampling). Table 2 lists the specification of all layers.

108  
109**Table 2:** FCN specification162  
163

	Name	Type	Input	Spec (size/stride)	Size CVPPP/MS-COCO	Size KITTI	Size Cityscapes	164
110	input	input	-	-	224 × 224 × 3	128 × 448 × 3	256 × 512 × 3	165
111	conv1-1	conv	input	3 × 3 × 3 × 32	224 × 224 × 32	128 × 448 × 64	256 × 512 × 64	166
112	conv1-2	conv	conv1-1	3 × 3 × 32 × 64	224 × 224 × 64	128 × 448 × 32	256 × 512 × 32	167
113	pool1	pool	conv1-2	max 2 × 2	112 × 112 × 64	64 × 224 × 64	128 × 256 × 64	168
114	conv2-1	conv	pool1	3 × 3 × 64 × 64	112 × 112 × 64	64 × 224 × 64	128 × 256 × 64	169
115	conv2-2	conv	conv2-1	3 × 3 × 64 × 96	112 × 112 × 96	64 × 224 × 96	128 × 256 × 96	170
116	pool2	pool	conv2-2	max 2 × 2	56 × 56 × 96	32 × 112 × 96	64 × 128 × 96	171
117	conv3-1	conv	pool2	3 × 3 × 96 × 96	56 × 56 × 96	32 × 112 × 96	64 × 128 × 96	172
118	conv3-2	conv	conv3-1	3 × 3 × 96 × 128	56 × 56 × 128	32 × 112 × 128	64 × 128 × 128	173
119	pool3	pool	conv3-2	max 2 × 2	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	174
120	conv4-1	conv	pool3	3 × 3 × 128 × 128	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	175
121	conv4-2	conv	conv4-1	3 × 3 × 128 × 128	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	176
122	conv4-3	conv	conv4-2	3 × 3 × 128 × 128	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	177
123	conv4-4	conv	conv4-3	3 × 3 × 128 × 128	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	178
124	conv4-5	conv	conv4-4	3 × 3 × 128 × 128	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	179
125	conv4-6	conv	conv4-5	3 × 3 × 128 × 128	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	180
126	conv4-7	conv	conv4-6	3 × 3 × 128 × 128	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	181
127	conv4-8	conv	conv4-7	3 × 3 × 128 × 256	28 × 28 × 256	16 × 56 × 256	32 × 64 × 256	182
128	pool4	pool	conv4-8	max 2 × 2	14 × 14 × 256	8 × 28 × 256	16 × 32 × 256	183
129	conv5-1	conv	pool4	3 × 3 × 256 × 256	14 × 14 × 256	8 × 28 × 256	16 × 32 × 256	184
130	conv5-2	conv	conv5-1	3 × 3 × 256 × 256	14 × 14 × 256	8 × 28 × 256	16 × 32 × 256	185
131	conv5-3	conv	conv5-2	3 × 3 × 256 × 256	14 × 14 × 256	8 × 28 × 256	16 × 32 × 256	186
132	conv5-4	conv	conv5-3	3 × 3 × 256 × 512	14 × 14 × 512	8 × 28 × 512	16 × 32 × 512	187
133	pool5	pool	conv5-4	max 2 × 2	7 × 7 × 512	4 × 14 × 512	8 × 16 × 512	188
134	deconv6-1	deconv	pool5	3 × 3 × 256 × 512/2	14 × 14 × 256	8 × 28 × 256	16 × 32 × 256	189
135	deconv6-2	deconv	deconv6-1 + conv5-3	3 × 3 × 256 × 512	14 × 14 × 256	8 × 28 × 256	16 × 32 × 256	190
136	deconv7-1	deconv	deconv6-2	3 × 3 × 128 × 256/2	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	191
137	deconv7-2	deconv	deconv7-1 + conv4-7	3 × 3 × 128 × 256	28 × 28 × 128	16 × 56 × 128	32 × 64 × 128	192
138	deconv8-1	deconv	deconv7-2	3 × 3 × 96 × 128/2	56 × 56 × 96	32 × 112 × 96	64 × 128 × 96	193
139	deconv8-2	deconv	deconv8-1 + conv3-1	3 × 3 × 96 × 192	56 × 56 × 96	32 × 112 × 96	64 × 128 × 96	194
140	deconv9-1	deconv	deconv8-2	3 × 3 × 64 × 96/2	112 × 112 × 64	64 × 224 × 64	128 × 256 × 64	195
141	deconv9-2	deconv	deconv9-1	3 × 3 × 64 × 64	112 × 112 × 64	64 × 224 × 64	128 × 256 × 64	196
142	deconv10-1	deconv	deconv9-2	3 × 3 × 32 × 64/2	224 × 224 × 32	128 × 448 × 32	256 × 512 × 32	197
143	deconv10-2	deconv	deconv10-1	3 × 3 × 32 × 32	224 × 224 × 32	128 × 448 × 32	256 × 512 × 32	198
144	deconv10-3	deconv	deconv10-2 + input	3 × 3 × 9 × 35	224 × 224 × 9	128 × 448 × 9	256 × 512 × 9	199
145								200

**Table 3:** External memory specification

196

144  
145  
146  
147  
148

	Name	Filter spec	Size CVPPP/MS-COCO	Size KITTI	Size Cityscapes	198
144	ConvLSTM	3 × 3	224 × 224 × 9	128 × 448 × 9	256 × 512 × 9	199
145						200

197  
198  
199

## D.2. External memory

201

## D.3. Box network

202

The box network takes in 9 channels of input directly from the output of the FCN. It goes through a CNN structure again and uses the attention vector predicted by the LSTM to perform dynamic pooling in the last layer. The CNN hyperparameters are listed in Table 4 and the LSTM and glimpse MLP hyperparameters are listed in Table 5. The glimpse MLP takes input from the hidden state of the LSTM and outputs a vector of normalized weighting over all the box CNN feature map spatial grids.

## D.4. Segmentation network

203

The segmentation networks takes in a patch of size 48 × 48 with multiple channels. The first three channels are the original image R, G, B channels. Then there are 8 channels of orientation angles, and then 1 channel of foreground heat map, all predicted by FCN. Full details are listed in Table 6. Constant  $\beta$  is chosen to be 5.

204  
205  
206  
207  
208  
209  
210

## References

211  
212  
213  
214  
215

- [1] P. Chattopadhyay, R. Vedantam, R. S. Ramprasaath, D. Batra, and D. Parikh. Counting everyday objects in everyday scenes. *CoRR*, abs/1604.03505, 2016. 1

216

**Table 4:** Box network CNN specification

217

Name	Type	Input	Spec (size/stride)	Size CVPPP/MS-COCO	Size KITTI	Size Cityscapes
input	input	-	-	224 × 224 × 9	128 × 448 × 9	256 × 512 × 9
conv1-1	conv	input	3 × 3 × 9 × 16	224 × 224 × 16	128 × 448 × 16	256 × 512 × 16
pool1	pool	conv1-2	max 2 × 2	112 × 112 × 16	64 × 224 × 16	128 × 256 × 16
conv1-2	conv	conv1-1	3 × 3 × 16 × 16	112 × 112 × 16	64 × 224 × 16	128 × 256 × 16
pool1	pool	conv1-2	max 2 × 2	56 × 56 × 16	32 × 112 × 16	64 × 128 × 16
conv2-1	conv	pool1	3 × 3 × 16 × 32	56 × 56 × 32	32 × 112 × 32	64 × 128 × 32
conv2-2	conv	conv2-1	3 × 3 × 32 × 32	56 × 56 × 32	32 × 112 × 32	64 × 128 × 32
pool2	pool	conv2-2	max 2 × 2	28 × 28 × 32	16 × 56 × 32	32 × 64 × 32
conv3-1	conv	pool2	3 × 3 × 32 × 64	28 × 28 × 64	16 × 56 × 64	32 × 64 × 64
conv3-2	conv	conv3-1	3 × 3 × 64 × 64	28 × 28 × 64	16 × 56 × 64	32 × 64 × 64
pool3	pool	conv3-2	max 2 × 2	14 × 14 × 64	8 × 28 × 64	16 × 32 × 64
conv3-1	conv	pool2	3 × 3 × 64 × 64	14 × 14 × 64	8 × 28 × 64	16 × 32 × 64
conv3-2	conv	conv3-1	3 × 3 × 64 × 64	14 × 14 × 64	8 × 28 × 64	16 × 32 × 64
pool3	pool	conv3-2	max 2 × 2	7 × 7 × 64	4 × 14 × 64	8 × 16 × 64

232

**Table 5:** Box network LSTM specification

233

Name	Size CVPPP/MS-COCO	Size KITTI	Size Cityscapes
LSTM	256	256	256
GlimpseMLP1	256	256	256
GlimpseMLP2	7 × 7	4 × 14	8 × 16

239

- [2] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 1

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324					378
325					379
326					380
327					381
328					382
329					383
330					384
331					385
332					386
333					387
334					388
335					389
336					390
337					391
338					392
339					393
340					394
341					395
342					396
343					397
344	Name	Type	Input	Spec (size/stride)	Size
345	input	input	-	-	48 × 48 × 13
346	conv1-1	conv	input	3 × 3 × 13 × 16	48 × 48 × 16
347	conv1-2	conv	conv1-1	3 × 3 × 16 × 32	48 × 48 × 32
348	pool1	pool	conv1-2	max 2 × 2	24 × 24 × 32
349	conv2-1	conv	pool1	3 × 3 × 32 × 32	24 × 24 × 32
350	conv2-2	conv	conv2-1	3 × 3 × 32 × 64	24 × 24 × 64
351	pool3	pool	conv2-2	max 2 × 2	12 × 12 × 64
352	conv3-1	conv	pool2	3 × 3 × 64 × 64	12 × 12 × 64
353	conv3-2	conv	conv3-1	3 × 3 × 64 × 96	12 × 12 × 96
354	pool3	pool	conv3-2	max 2 × 2	6 × 6 × 96
355	deconv4-1	deconv	pool3	3 × 3 × 64 × 96/2	12 × 12 × 64
356	deconv4-2	deconv	deconv4-1 + conv3-1	3 × 3 × 64 × 128	12 × 12 × 64
357	deconv5-1	deconv	deconv4-2 + conv2-2	3 × 3 × 32 × 128/2	24 × 24 × 32
358	deconv5-2	deconv	deconv5-1 + conv2-1	3 × 3 × 32 × 64	24 × 24 × 32
359	deconv6-1	deconv	deconv5-2 + conv1-2	3 × 3 × 16 × 64/2	48 × 48 × 16
360	deconv6-2	deconv	deconv6-1 + conv1-1	3 × 3 × 16 × 32	48 × 48 × 16
361	deconv6-3	deconv	deconv6-2 + input	3 × 3 × 1 × 29	48 × 48 × 1
362					414
363					415
364					416
365					417
366					418
367					419
368					420
369					421
370					422
371					423
372					424
373					425
374					426
375					427
376					428
377					429
					430
					431