

看图说话机器人-densecap

Dense caption-经典论文

[1]DenseCap: Fully Convolutional Localization Networks for Dense Captioning: *Justin Johnson, Andrej Karpathy, Li Fei-Fei*

[2]Dense Captioning with Joint Inference and Visual Context: *Linjie Yang Kevin Tang Jianchao Yang Li-Jia Li*



Dense caption任务是image caption和object detection任务相结合。

其中object detection任务一般采用的是Faster R-CNN 网络, Caption Model则采用single LSTM 或 multil LSTM进行构建, 这两个部分间也有多种结合的方式。

Dense caption[1]

densecap模型的结构如图,
网络由4部分构成:

- CNN 特征提取网络
- Localization Layer
- Recognition Network
- RNN caption Model

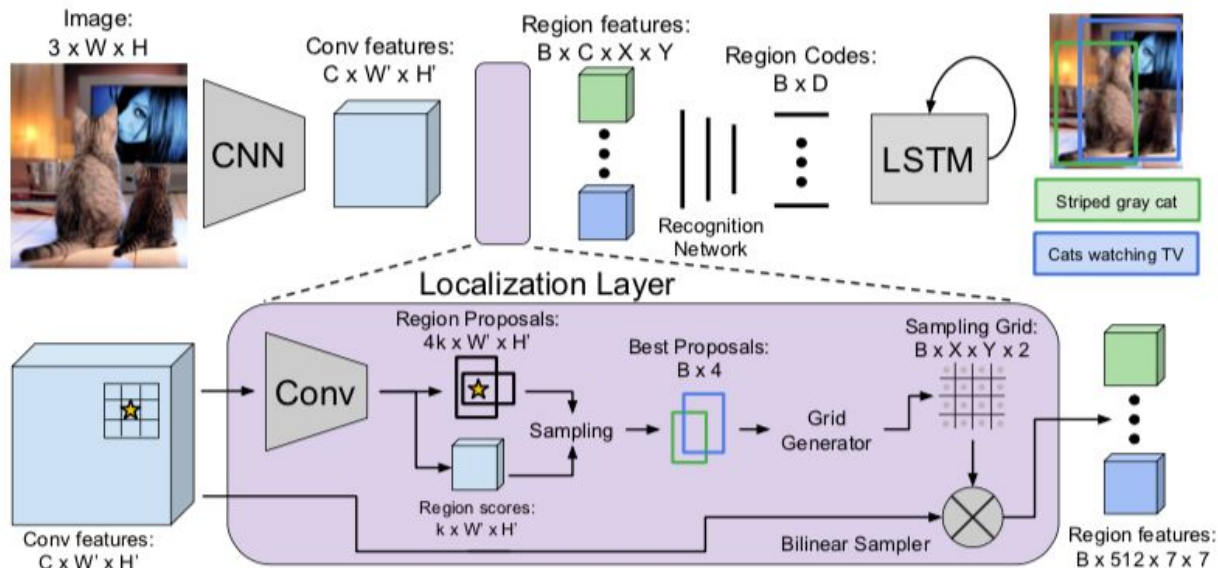
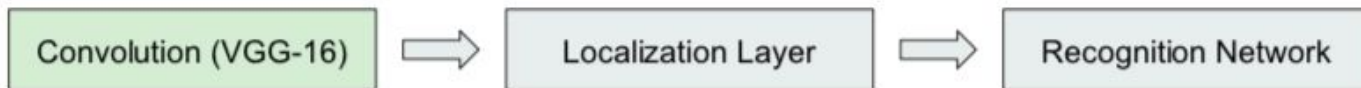


Figure 2. Model overview. An input image is first processed a CNN. The Localization Layer proposes regions and smoothly extracts a batch of corresponding activations using bilinear interpolation. These regions are processed with a fully-connected recognition network and described with an RNN language model. The model is trained end-to-end with gradient descent.

Dense caption[1]

图像特征提取网络论文中选取的是VGG-16,对输入图片提取出其特征向量, 输入到**Localization layer**



Convolution Network

- VGG-16 network for state-of-the-art performance in converting input images to tensor feature matrix.
- 13 layers of 3×3 convolutions interspersed with 5 layers of 2×2 max pooling.
- Input of $3 \times W \times H$ with output tensor features of $512 \times W/16 \times H/16$.

Dense caption[1]

Localization layer 能够识别感兴趣的区域并且从每一个区域平滑得提取一个固定尺寸的图像向量表示。本文基于 Faster R-CNN但是用双线性插值代替了 RoI pooling 机制。这一变化的好处在于打开了预测仿射或者变形候选区域的可能性而不仅仅是bounding box。

输入 => 来自VGG16卷积网络的特征图谱 $C \times W' \times H' \times C \times W' \times H'$ (任意size)

输出 => 输出B个候选区域的表征向量(定长), 在每个特征向量都包含下面三个关键信息:

1. 候选区域的坐标:输出形式是一个 $B \times 4 \times B \times 4$ 的矩阵, 每行代表一个候选区域的坐标
2. 候选区域的置信分数:一个长度为 $B \times B$ 的一维列向量, 向量内每个元素都给出了候选区域的得分。得分越高说明越可能是真实区域
3. 候选区域的特征:输出形式为 $B \times C \times X \times Y \times B \times C \times X \times Y$ 的特征集合, 这里B代表区域个数, $X \times Y \times X \times Y$ 表示特征图谱的大小(注意, 这里的size已经是固定的), CC代表特征的维度

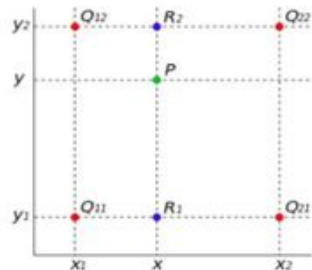
Dense caption[1]

Localization layer 中的Faster R-CNN 用双线性插值代替了原有的RoI pooling机制, 但在结构中还是称其为pooling层

Bilinear interpolation

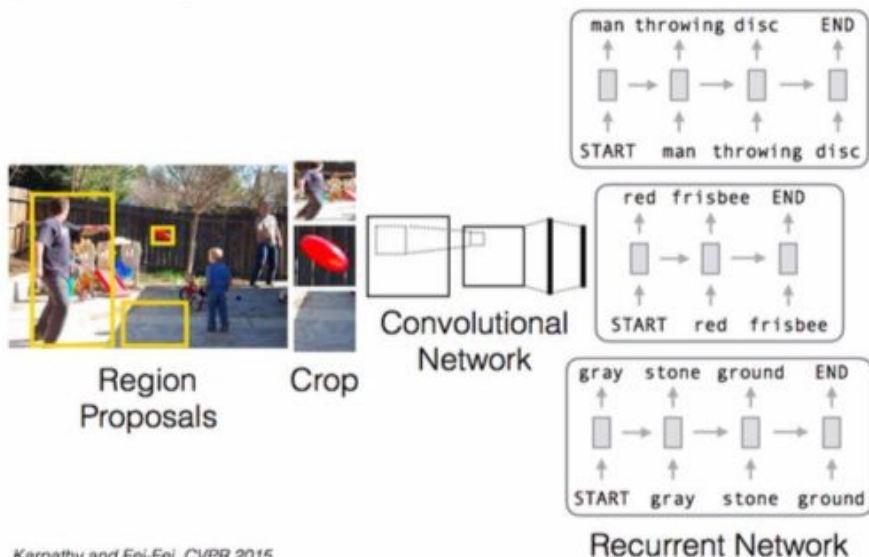
- Allows for backpropagation of error to previous layers.

$$f(x, y) = \frac{1}{(x_2 - x_1)(y_2 - y_1)} (f(Q_{11})(x_2 - x)(y_2 - y) + f(Q_{21})(x - x_1)(y_2 - y) + f(Q_{12})(x_2 - x)(y - y_1) + f(Q_{22})(x - x_1)(y - y_1))$$



Dense caption[1]

Recognition network



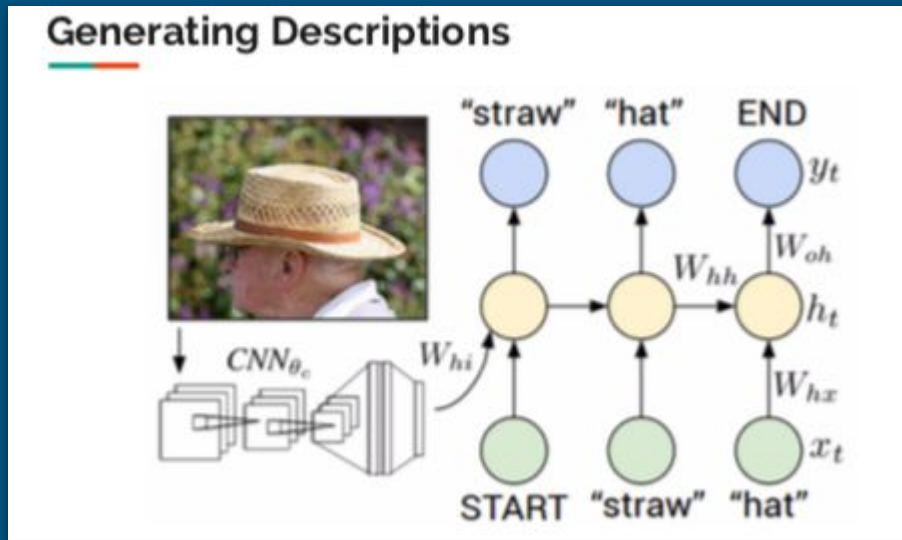
Karpathy and Fei-Fei, CVPR 2015

Recognition Network网络的作用就是将来自定位层的每个region的特征falten为一个向量并将其传至两个全连接层，这两层都是使用了ReLU单元并且使用了Dropout。每一个region从而能够产生一个维数为 $D = 4096$ 的code，里面包含了这一个region中的visual appearance。所有positive region的codes被收集到了一个大小为 $B \times D$ 的矩阵中，然后被传送至RNN语言模型。

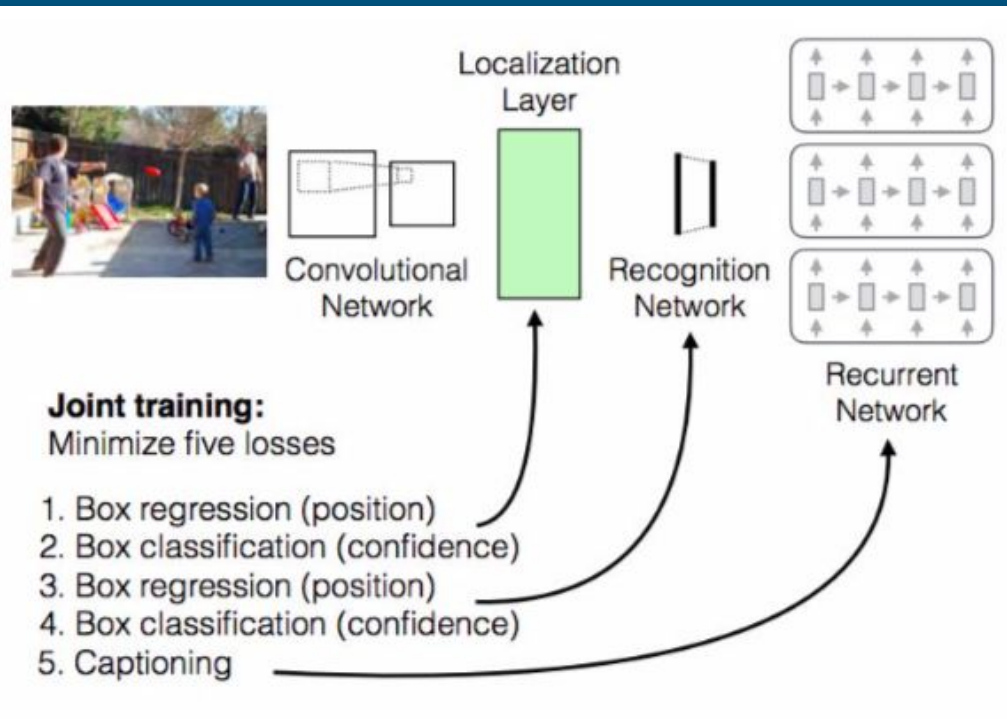
另外，Recognition Network也可以再次完善每个region proposal的置信度和位置。这两者来自于每个region所对应的 D 维code的线性变换。

Dense caption[1]

将得到的图像区域特征向量输入到我们的RNN模型中并生成描述 这个过程和我们在show and tell的im2txt模型中是一样的



Dense caption[1]



整个网络模型的Loss function:

采样后的positive和negative regions里, 对于置信度, 文中使用的是binary logistic loss(被使用两次);

于box regression, 使用的是smooth L1 loss来变换坐标空间(被使用两次);

语言模型中每个时刻用的是 cross-entropy。

有的损失函数都是通过RNN中的batch size和sequence length来正则化的。

Dense caption[1]

(binary logistic loss)

$$l(w, b) = - \sum_{i=1}^m \ln P(y_i | x_i; w, b)$$

$$P(y = 1 | x) = \frac{e^{w^T x + b}}{1 + e^{w^T x + b}}$$

$$P(y = 0 | x) = \frac{1}{1 + e^{w^T x + b}}$$

(cross-entropy)

$$L(I, S) = - \frac{1}{T} \sum_{t=1}^T \log p_t(x_t)$$

(smooth L1 loss)

$$L_{loc}(t^u, v) = \sum_{i \in \{x, y, w, h\}} \text{smooth}_{L_1}(t_i^u - v_i)$$

$$\text{smooth}_{L_1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

各loss的计算公式

Dense caption[1]

对于densecap的训练: 用预训练ImageNet的权重来初始化CNN, 其他权重的取值来源于标准差为0.01的高斯分布。CNN中权重的训练使用了momentum为0.9的SGD, 其他的权重使用了Adam算法。学习速率分别为 $\beta_1=0.9$ 和 $\beta_2=0.99$ 。CNN的fine-tuning在迭代一轮后就开始了, 但是不包含CNN中的前四个卷积层的微调。

Results (Image Retrieval)

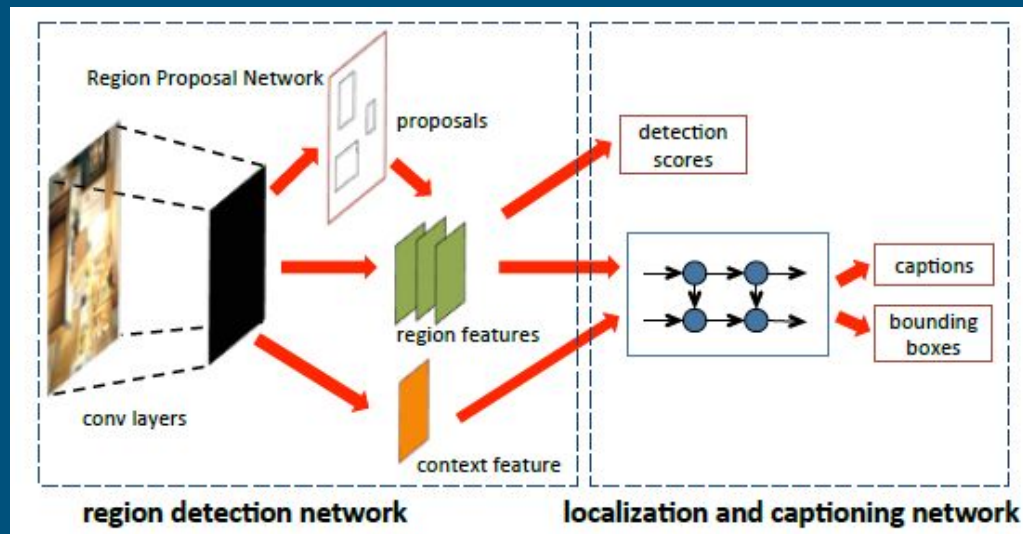
	Ranking				Localization			
	R@1	R@5	R@10	Med. rank	IoU@0.1	IoU@0.3	IoU@0.5	Med. IoU
Full Image RNN [21]	0.10	0.30	0.43	13	-	-	-	-
EB + Full Image RNN [21]	0.11	0.40	0.55	9	0.348	0.156	0.053	0.020
Region RNN [21]	0.18	0.43	0.59	7	0.460	0.273	0.108	0.077
Our model (FCLN)	0.27	0.53	0.67	5	0.560	0.345	0.153	0.137

- Performs much better than most other models when doing natural language queries to retrieve images.

Dense caption[2]

Dense Captioning with Joint Inference and Visual Context: Linjie Yang Kevin Tang Jianchao Yang Li-Jia Li

该论文模型结构如下图：



模型由两个部分组成：
a region detection network 和
a localization and captioning
network。
前者为采用Faster-RCNN的区域监
测网络，后者进行关键区域定位和描
述生成任务。

Dense caption[2]

Dense Captioning with Joint Inference and Visual Context: Linjie Yang Kevin Tang Jianchao Yang Li-Jia L

这篇论文是从李飞飞团队的densecap出发, 在其基础上进行创新, 论文提出了densecap任务的两个难点:

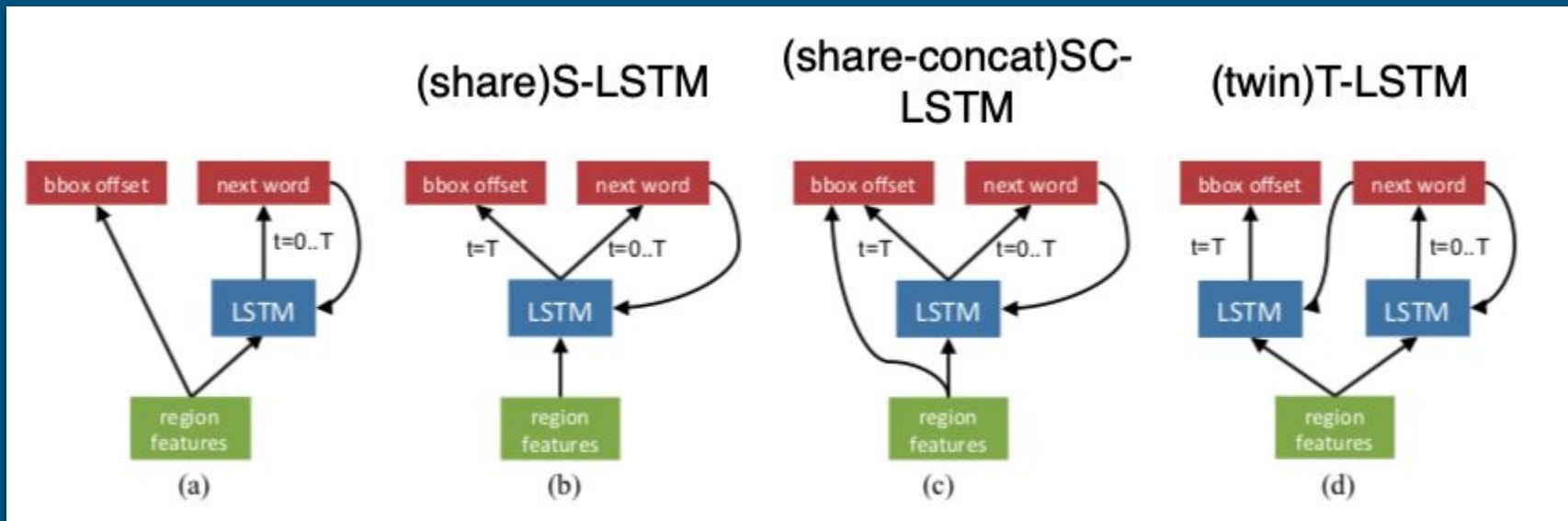
- 1、object detection获得的bounding box往往会发生堆叠重合的情况, 因此要根据描述来准确定位区域;
- 2、很多区域被描述, 由于有多种描述, 导致部分区域的内容变得模糊不清。

本论文中提出两个策略:**joint inference** 和 **context fusion**来对以上问题进行改进。以上描述的2种模式在效果上, 同encoder-decoder结构中的attention机制有一定程度的相似, 但计算方式是完全不同的, 接下来我们对这2个方法进一步的了解。

Dense caption[2]

joint inference:

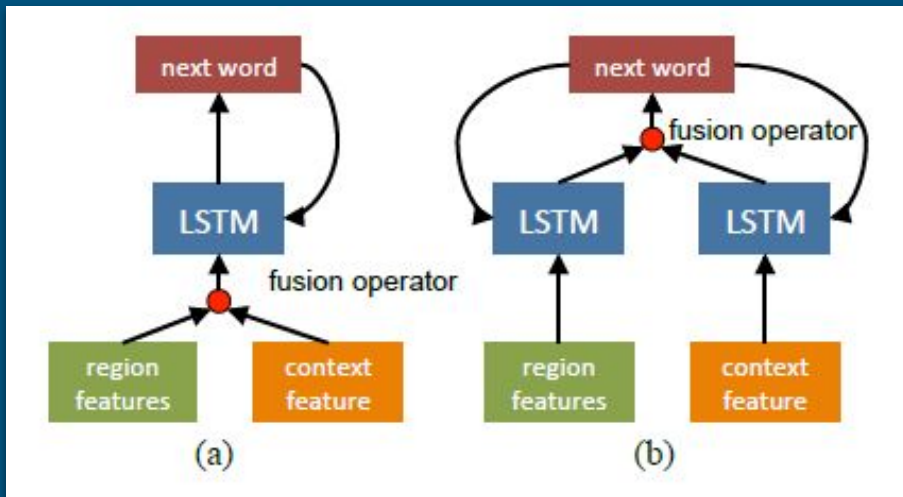
为了让bounding box的定位受相关联图片区域的描述的影响, 提高定位的准确性, 我们利用LSTM不仅用来生成描述, 还用来对bounding box的定位进行修饰, 论文列出以下很多种方式, 如下所示:



im2txt 进阶 Dense captioning[2]

context fusion:

为了将context feature和local feature结合起来, 论文选用RoI pooling feature作为图片的visual context feature较好的形式, 但是context feature并没有直接和bounding box直接关联, 仅仅使用context feature协助生成caption, 再利用上面的joint inference来影响local feature。本文为了实现此功能, 有如下两种形式ealy-fusion和late-fusion:



im2txt 进阶 Dense captioning[2]

有了以上两种创新结构，我们可以有多种组合形式。最后模型的损失函数如下：

$$L = L_{cap} + \alpha L_{det} + \beta L_{bbox}$$

L-caption表示caption model的交叉熵损失函数，计算的是预测下一个单词的出现概率的损失

L-det表示detection的交叉熵损失函数，计算的是前景和背景的二分类交叉熵损失

L-bbox 表示bounding box的L1正则损失函数

α 、 β 为加权系数

Dense caption[2]

Table 3: The mAP performance of integrated models with combinations of joint inference models and context fusion structures on Visual Genome V1.0.

model		S-LSTM	SC-LSTM	T-LSTM
early-fusion	$[\cdot, \cdot]$	6.74	7.18	8.24
	\oplus	6.54	7.29	8.16
	\otimes	6.69	7.04	8.19
late-fusion	$[\cdot, \cdot]$	7.50	7.72	8.49
	\oplus	7.19	7.47	8.53
	\otimes	7.57	7.64	8.60

Table 4: The mAP performance of different dense captioning models on Visual Genome V1.2.

model		baseline	S-LSTM	T-LSTM
no context			6.44	8.16
late-fusion	$[\cdot, \cdot]$	6.98	7.76	9.03
	\oplus		7.06	8.71
	\otimes		7.63	8.52

同时论文中列出了网络结构中的joint inference和context fusion不同类型的结果对比, 可以看出T-LSTM和Late-fusion是性能较高的一组。

所以在我们的网络结构中实现的是T-LSTM和Late-fusion, 其中Late-fusion采用的sum的方式进行实现的。

在和densecap论文的模型的性能对比上, 本论文实现的效果达到了73%的提高。

	Visual Genome V1.0			V1.2
Model	Johnson <i>et al.</i> [20]	Ours	Gain	Ours
mAP	5.39	9.31	73%	9.96

Dense caption工程实现

❖ 代码参考

- The Faster-RCNN framework: [tf-faster-rcnn](#) by [endernewton](#)
- Py-faster-rcnn framework: [py-faster-rcnn](#) by [Ross Girshick](#)
- Dense Captioning with Joint Inference and Visual Context's code: [densecap](#)
- Dense Cap official code: [densecap](#) by [Justin Johnson](#)
- [Densecap](#) by InnerPeace-Wu

我们的代码最终实现的模型为论文 Dense Captioning with Joint Inference and Visual Context 中提出的模型。Joint Inference部分实现的是T-LSTM结构context fusion则对late fusion中的sum方式进行的实现。

Dense caption工程实现

代码实现过程中, 对于模型最重要的部分 Faster R-CNN的代码我们选用了github上star较多的采用了tensorflow进行实现的仓库, 并且在做的过程中也发现有部分代码需要从py-faster-rcnn仓库中进行引用, 这部分代码主要是对RPN结构进行处理的部分。论文[2]的官方实现是caffe下实现的代码, 其中模型部分的代码需要我们在tensorflow下重构。对于caption model在进行test or inference时, 要采取beam search对生成的语句进行筛选, 这里引用了im2txt中的代码。

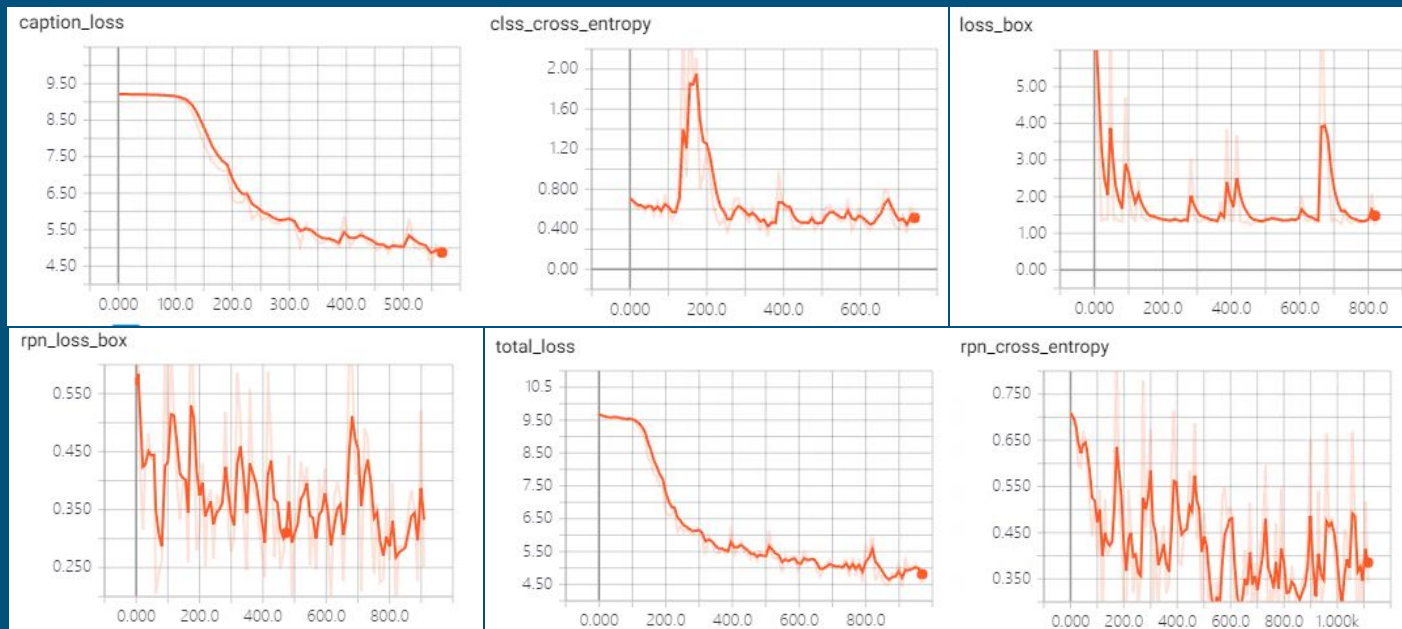
由于实现的模型由几个部分组成, 而这几个部分在网络训练过程中要进行调参, 所以模型调参的参数是比较多的, 所以必须采取统一的配置文件对参数进行设置, 训练过程中对该配置文件进行调整来适应不同阶段的调参需求。参数虽然较多, 但大部分可以不用调整, 保持官方代码中对应的值。

由于引用了ms coco的evaluate模块(pycocoevalcap), 在进行val的时候需要python2.7的环境, train和inference则可以跑在任意版本下。

由于我们的运算设备有足够的内存和GPU, 代码中设置数据集batchsize为256, 在一些内存有限制的设备上则需要调小。

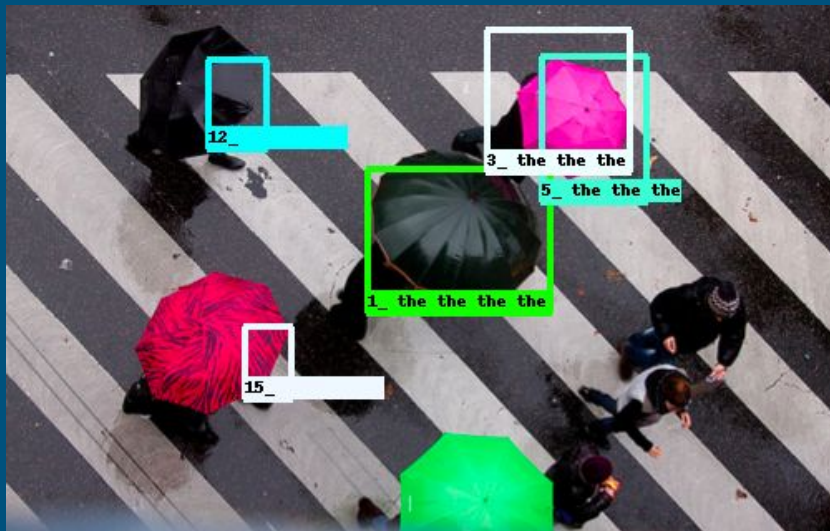
Dense caption训练调参

训练的过程分几个阶段进行, 开始阶段会引入Faster R-CNN的预训练模型, 第一阶段不对detection结构的参数进行调整, 调整其他结构的参数, 训练的收敛非常快(初始学习率为0.001):



Dense caption训练调参

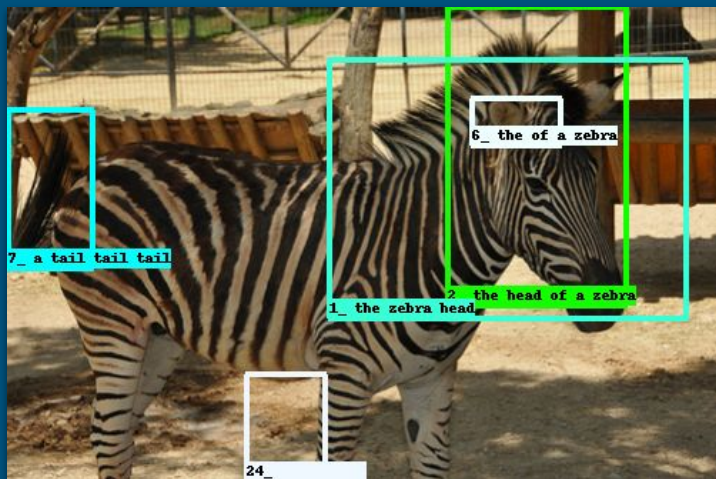
在第一个阶段模型各项loss收敛至平稳后，则开始对detection model的参数进行训练调整，同时需要把学习率调 $\frac{1}{4}$ ，继续训练至loss收敛。前2个阶段中我们用tensorboard可视化训练中的图像可以看到，模型在detection部分的提升是较为明显的。



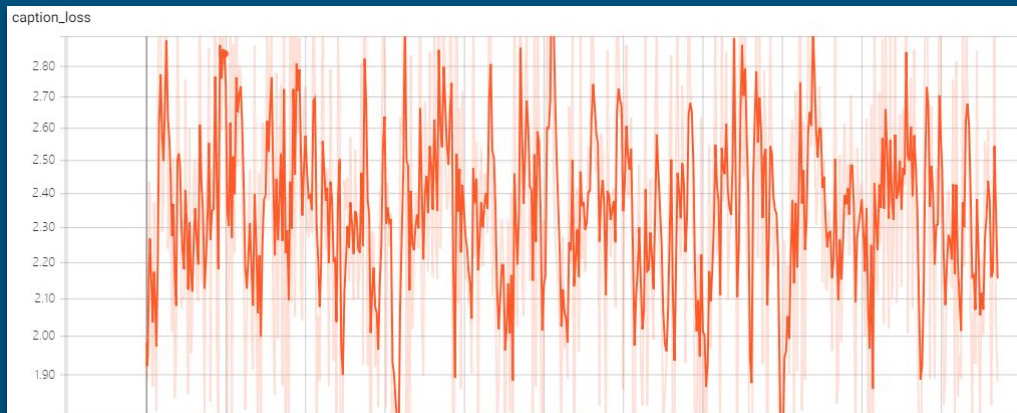
Dense caption训练调参

在对detection model进行调参至收敛后，则开始对caption model进行调参数，这里调整的设置主要是context fusion的开启，并锁定 cnn网络的参数，学习率继续衰减 $\frac{1}{4}$ ，同样训练至loss收敛。

最后阶段的训练则在上阶段的基础上取消对cnn网络参数的锁定，学习率继续衰减 $\frac{1}{4}$ ，开启训练。



网络50K迭代训练captionloss如图，值在2.4起伏不定，没有收敛



Dense caption成果

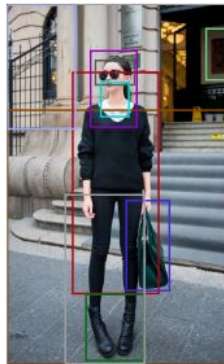
由于模型中由多个不同的结构构成, 所以模型的参数较多, 并且在调参过程中分了几个不同的阶段对参数进行调试, 较im2txt调参过程更加复杂。(右图为训练50K的模型效果。)

另外目前虽然loss没有达到较小的收敛区间, 一直在2.X左右, 但网络效果已经呈现不错的性能, MAP的得分为8.25。

生成的句子比较丰富, 但是 we 希望能生成比较流畅完整的句子。



a man wearing a blue shirt. a woman wearing sunglasses. a pair of sunglasses. the head of a man. the man is wearing a black shirt. a white shirt on a woman. a woman wearing sunglasses. a yellow sign on the wall. the hand of a person. the shirt is black. a man in a black shirt. a black chair in the background. the woman has blonde hair. a building in the background. the hand of a man. people in the background. the arm of a man.



a woman wearing a black jacket. the bag is green. black sunglasses on a woman. a pair of black boots. a woman walking on the sidewalk. a woman wearing sunglasses. the woman is wearing black pants. picture on the wall. black metal railing on the side of the building. the woman is wearing a necklace.

实验效果

Dense caption进阶

[A Hierarchical Approach for Generating Descriptive Image Paragraphs: Justin Johnson, Andrei Karpathy, Li Fei-Fei](#)

这是李飞飞团队对densecap领域提出的第二篇论文，稳重提出的模型将图像作为输入，生成描述它的自然语言段落，并且旨在利用图像和Paragraphs的组合结构，起结构如下。

Region Detector: 我们首先通过目标检测网络对输入图像提取出图像中的propose regions of interest区域向量。

Region Pooling: 然后在这些区域向量中聚合特征(max pooling)以产生富有表达图像语义的合并表示。(其中论文提出了相关同期一些工作，表示max pooling不会丢失需要表达的重要信息)

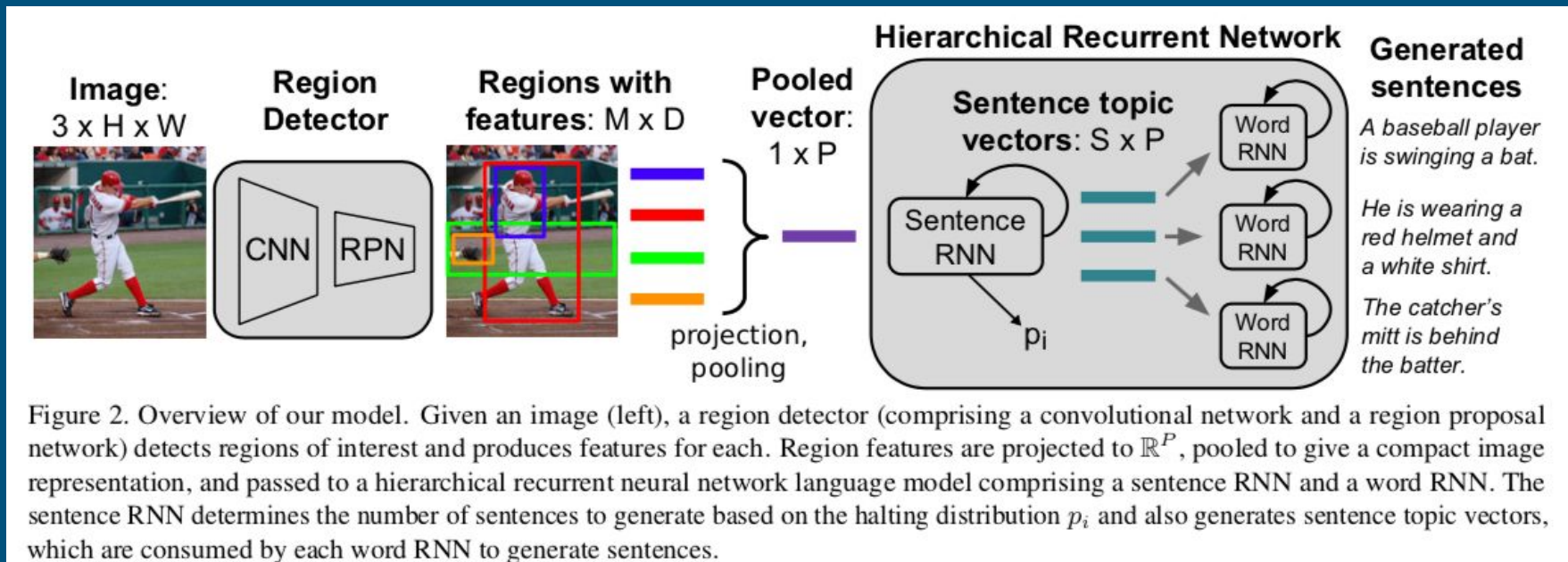
Hierarchical Recurrent Network，该结构包含Sentence RNN和Word RNN

Sentence RNN: 每个timestep都接收池化区域特征向量作为输入，其初始状态值为0，该结构决定在结果段落中生成多少句子，并为每个句子产生输入主题向量。

Word RNN: 根据输入的句子主题向量，生成每个句子的单词。最终生成一个段落描述。

Dense caption进阶

论文中提出的模型结构如图，该网络也为我们的模型改进提供了方向和思路。



Dense caption总结

A Hierarchical Approach for Generating Descriptive Image Paragraphs: Justin Johnson, Andrej Karpathy, Li

Fei-Fei 这篇论文为我们提供了改进的方向, 同时论文中也给出了李飞飞团队对该任务制作的数据集, 他们在visual genome基础上补充了图片的段落描述标签, 为我们进行下阶段的任务提供了保障, 对于我们当前的模型, 也有很多需要 进行优化的地方比如loss收敛, 生成重复句子, 生成句子的 时间较长导致的服务响应较慢。