



华南师范大学
SOUTH CHINA NORMAL UNIVERSITY

2019—2020年度学生课外科研一般课题

结题报告

课题名称：基于深度学习的图像描述

第一作者：刘杰聪

课题组成员：李海宏、谢淦煜、赖乐贤

指导教师：梁军

所在院系：软件学院

结项日期：2020年4月

共青团华南师范大学委员会制

2020.4

华南师范大学 2019—2020 年度学生课外科研一般课题

结题申请表

课题名称		基于深度学习的图像描述				
课题形式		<input type="checkbox"/> 个人课题 <input checked="" type="checkbox"/> 集体课题				
学科类别		<input type="checkbox"/> 哲学 <input type="checkbox"/> 社会 <input type="checkbox"/> 法律 <input type="checkbox"/> 经济 <input type="checkbox"/> 管理 <input type="checkbox"/> 教育 <input type="checkbox"/> 政治 <input type="checkbox"/> 历史 <input type="checkbox"/> 文学 <input type="checkbox"/> 艺术 <input type="checkbox"/> 体育 <input type="checkbox"/> 数理 <input checked="" type="checkbox"/> 信息技术 <input type="checkbox"/> 机械与控制 <input type="checkbox"/> 生命科学 <input type="checkbox"/> 能源化工				
课题类别		<input type="checkbox"/> 哲学社会科学类调查报告和学术论文 <input type="checkbox"/> 自然科学类学术论文 <input checked="" type="checkbox"/> 科技发明制作类 <input type="checkbox"/> 创作成果				
指导教师		姓名	职称	学 院	联系电话	
		梁军	工程师	软件学院	13632368330	
作者简介	第一作者	姓名	刘杰聪		性别	男
		出生年月	1999 年 1 月		学历	本科
		学院 年级	软件学院 2018 级			
		联系地址	茂名市电白区观珠镇			
		联系电话	13580070189			
	其他作者	姓名	性别	学历	所在学院	
		李海宏	男	本科、2018 级	软件学院	
		谢淦煜	男	本科、2018 级	软件学院	

	者	赖乐贤	男	本科、2018 级	软件学院
结 题 申 请	<p>本小组已经完成院级课题《基于深度学习的图像描述》，申请结题</p> <p>课题负责人签字：刘杰聪</p> <p>2020 年 4 月 25 日</p>				
经 费 支 出 概 况	支出内容				金 额
	<u>《深度学习》*4</u>				288
	<u>《深度学习入门之 pytorch》*4</u>				56
	<u>《 Python 从入门到精通》*4</u>				170
	(备注：以上均为小组自费)				0
	总计				514
课 题 指 导 教 师 评 定 意 见	<p>(1) 评定内容说明：课题研究是否按计划开展，研究内容有无明显抄袭和 不实之处，研究成果质量如何，研究者是否主动征求导师指导意见并就研究相关问题与导师交流等。</p> <p>(2) 评定情况及分数说明：85 分以上为优秀，70—85 分为良好，60—70 为及格，60 以下为不及格。</p>				
	<p>综合评定：</p> <p>该项目主要识别图像，并用一个语句来描述该图像的含义。模型是基于 CNN-LSTM 的编码器-解码器框架，并利用注意力机制提取重要区域，并优化了界面，实验表明，该项目有较好的识别效果。该项目达</p>				

	<p>到结题要求，同意结题。</p> <p>综合得分：88</p> <p>指导教师签字：梁军</p> <p>2020 年 4 月 29 日</p>
学院 专家 评审 委员会 意见	<p>主审专家签字：</p> <p>年 月 日</p>
学院 课外 科技 创新 领导 小组 意见	<p>领导签字：</p> <p>部门盖章：</p> <p>年 月 日</p>
学校 课外	

科技创新领导小组意见	<div>领导签字：</div> <div>部门盖章：</div> <div>年 月 日</div>
备注	<div>1、本报告中指导老师评定意见、学院专家评审委员会意见和学院课外科技创新领导小组意见均须手写，不得用打印，否则该申请表作废；</div> <div>2、除 1 中所述三项内容外其余内容均须用计算机打印，文字力求精练、准确，不得使用非规范用语；</div> <div>3、课题形式、学科类别、课题类别的方框均采用涂黑（■）的形式进行选择</div> <div>4、提交本申请表时需同时提交本申请表的电子文档；</div> <div>5、结题成果请按《学生课外科研课题成果（文本）格式要求》另附材料说明。</div>

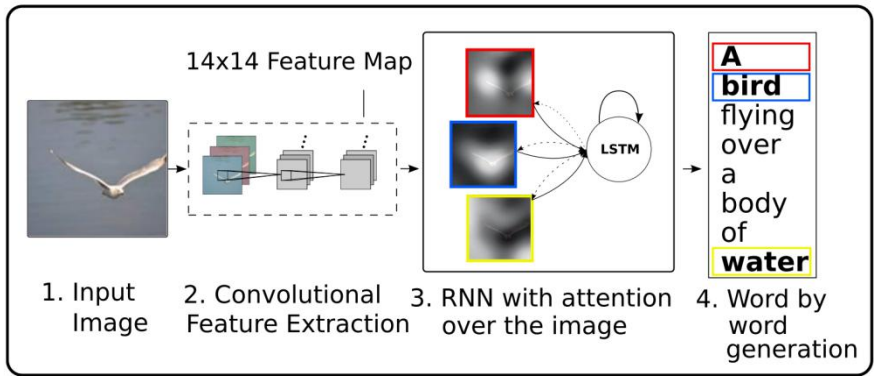
课题内容简介及课题成果

课题内容简介：

图像描述，顾名思义，就是给计算机提供图片，让计算机识别出图片，并以自然语言表述出来。基于深度学习的方法的优势是可以直接从大量的训练数据中学习图像到描述语句的映射，实现端到端的训练，并能产生更精确的图像描述，在性能上远远优于传统方法。

本小组采用的深度学习框架是 pytorch，其设计思路是线性、直观且易于使用的；采用的模型是基于 CNN-RNN（Encoder-Decoder）的编码器-解码器框架：对整个图像描述生成过程进行建模，分成两个部分：基于卷积神经网络的图像编码器和基于循环神经网络的句子解码器。编码，即将输入序列转化成一个固定长度的向量；解码，即将之前生成的固定向量再转化成输出序列。模型首先利用 CNN 对图像进行有效的特征提取与编码，然后利用 RNN 及其变体 LSTM 生成描述句子。另外本课题引入了注意力机制：解码器的每个时刻在生成一个单词的时候，注意力机制模块根据之前已经预测出的单词信息来提示当前时刻应该关注图像的哪些重要区域，而不是漫无目的地关注整幅图像，然后利用与本时刻相关性高的图像区域特征来生成单词。通俗地来讲，注意力机制就是让网络在解码的时候能够“集中注意力”在编码输出的某些部分上。注意力机制的引入大大提高了我们采用基于编码器-解码器框架的图像描述模型的性能。

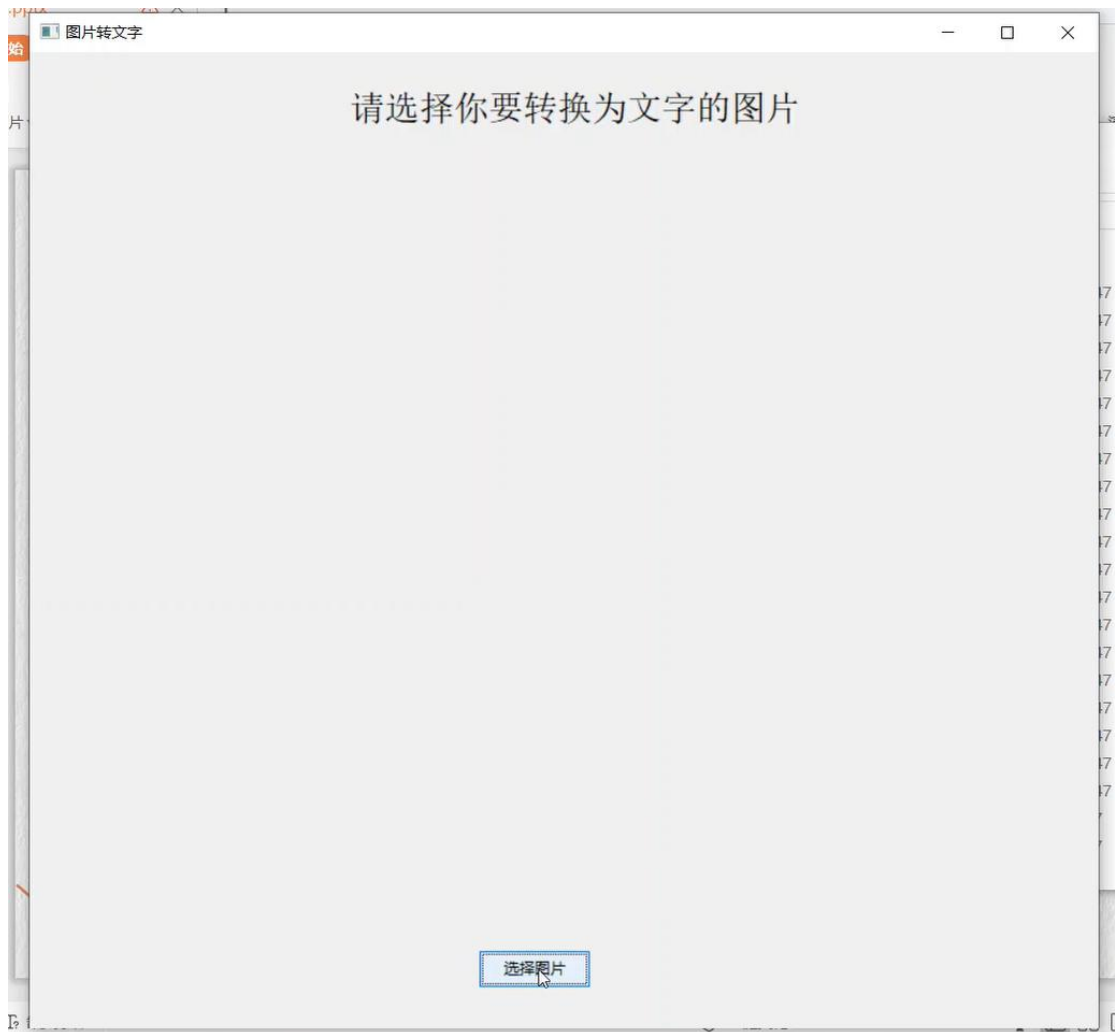
Encoder-Decoder 框架示意图：



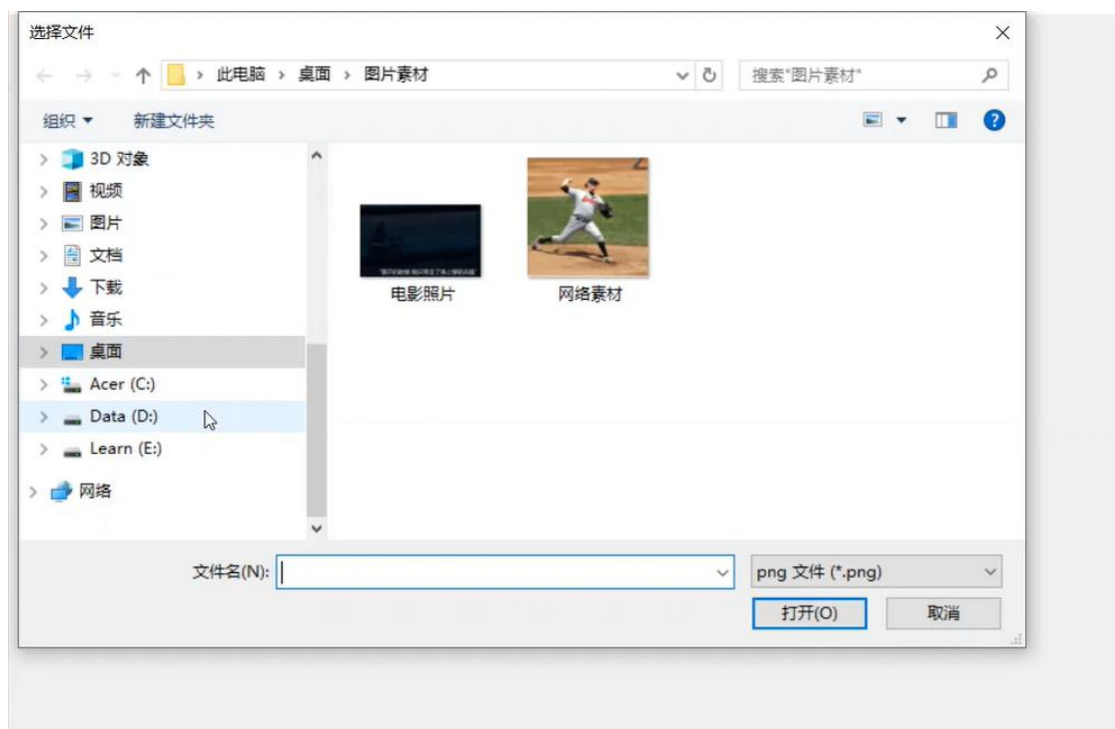
课题成果：

本小组课题成果是一个打包好的软件，实现了申报时的预期目标，由于代码量较大（接近 4G），所以就不展示代码了，下面是软件运行截图，详细的软件演示可见演示视频。

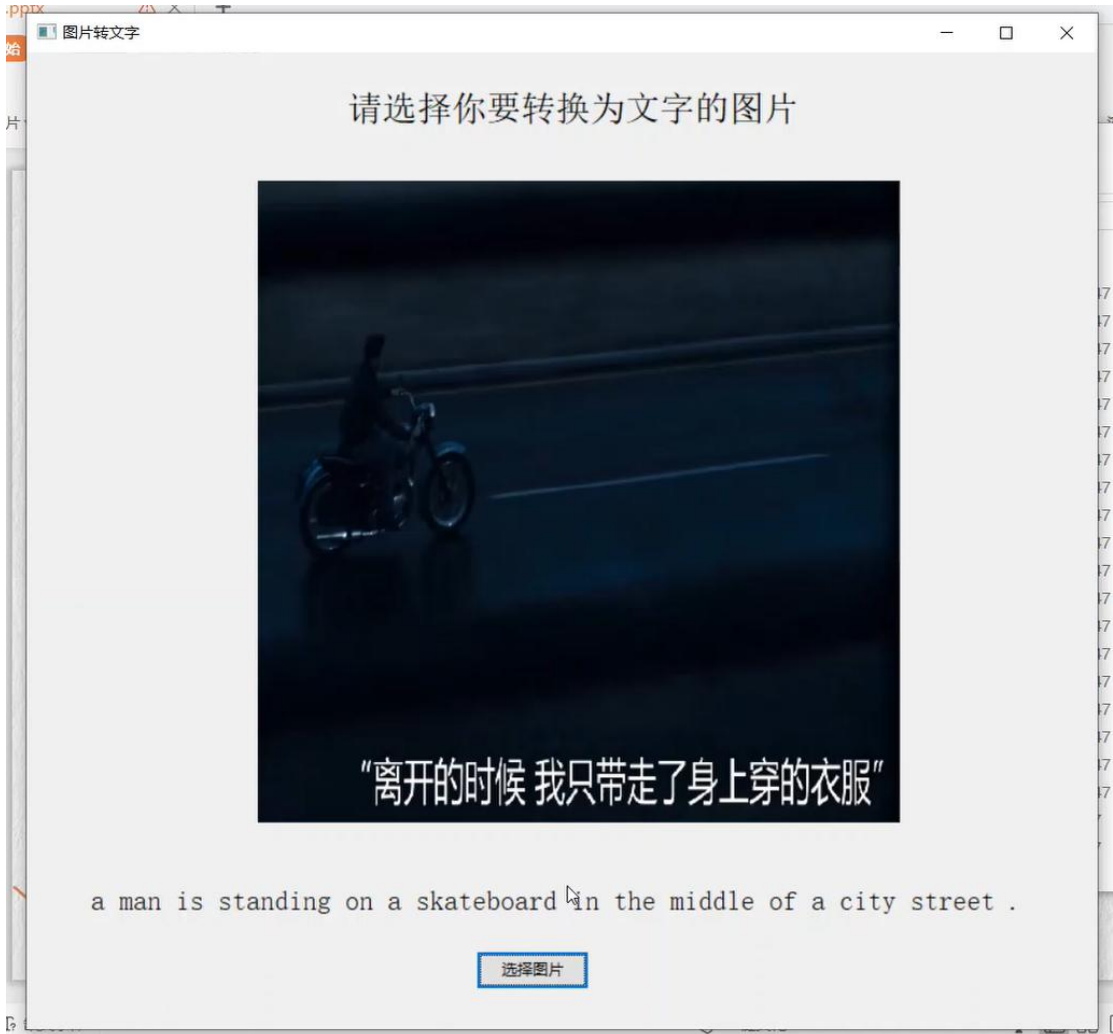
这是软件的打开界面，比较简洁：



打开后，点击选择图片，会弹出文件选择框供选择：



选择完成后，第一次加载会比较久，因为软件会调用一些模型文件，然后进行识别描述，将描述句子显示在图片下方：



图片的描述评估分数 BLEU 分数已达到 69，接近了参考的相关论文分数 71：

arXiv:1502.03044v3 [cs.LG] 19 Apr 2016

Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

Kehin Xu
Jimmy Lei Ba
Ryan Kiros
Kyunghyun Cho
Aurene Courville
Ruslan Salakhutdinov
Richard S. Zemel
Yoshua Bengio

KUANYI.XU@UNIONEDUCAL.CA
JIMMY@PSI.UTORONTO.CA
KIROS@PSI.UTORONTO.CA
KHYUNGHYUN.CHO@UNIONEDUCAL.CA
AURENE.COURVILLE@UNIONEDUCAL.CA
RUSLAN.SALAKHUTDINOV@UNIONEDUCAL.CA
ZEMEL@PSI.UTORONTO.CA
YBENGIO@THE.UBC.CA

Abstract

Inspired by recent work in machine translation and object detection, we introduce an attention based model that automatically learns to describe the content of images. We describe how we can train this model in a deterministic manner using standard backpropagation techniques and stochastically by maximizing a variational lower bound. We also show through visualization how the model is able to automatically learn to fix its gaze on salient objects while generating the corresponding words in the output sequence. We validate the use of attention with state-of-the-art performance on three benchmark datasets: Flickr8K, Flickr30K and MS COCO.

1. Introduction

Automatically generating captions of an image is a task very close to the heart of scene understanding—one of the primary goals of computer vision. Not only must caption generation models be powerful enough to solve the complex vision challenges of determining which objects are in an image, but they must also be capable of capturing and expressing their relationships in a natural language. For this reason, caption generation has long been viewed as a difficult problem. It is a very important challenge for machine learning algorithms, as it amounts to mimicking the remarkable human ability to compress huge amounts of salient visual information into descriptive language.

Despite the challenging nature of this task, there has been a recent surge of research interest in attacking the image caption generation problem. Aided by advances in training neural networks (Kreishutzky et al., 2012) and large classification datasets (Russakovsky et al., 2014), recent work

Figure 1. Our model learns a word/image alignment. The visualized attention maps (1) are explained in section 3.1 & 3.4.

1. Heat Map 2. Convolutional 3. Soft self-attention 4. Hard self-attention

has significantly improved the quality of caption generation using a combination of convolutional neural networks (convnets) to obtain *spatial representation* of images and recurrent neural networks to decode these representations into natural language sentences (see Sec. 2).

One of the most curious facts of the human visual system is the presence of attention (Rensink, 2000; Carbotte & Dhanasekaran, 2002). Rather than compress an entire image into a static representation, attention allows for salient features to dynamically come to the forefront as needed. This is especially important when there is a lot of clutter in an image. Using representations such as those from the top layers of a convnet that distill information in image down to the most salient objects is one effective solution that has been widely adopted in previous work. Unfortunately, this has one potential drawback of losing information which could be useful for richer, more descriptive captions. Using more low-level representations can help preserve this information. However working with these features necessitates a powerful mechanism to steer the model to information important to the task at hand.

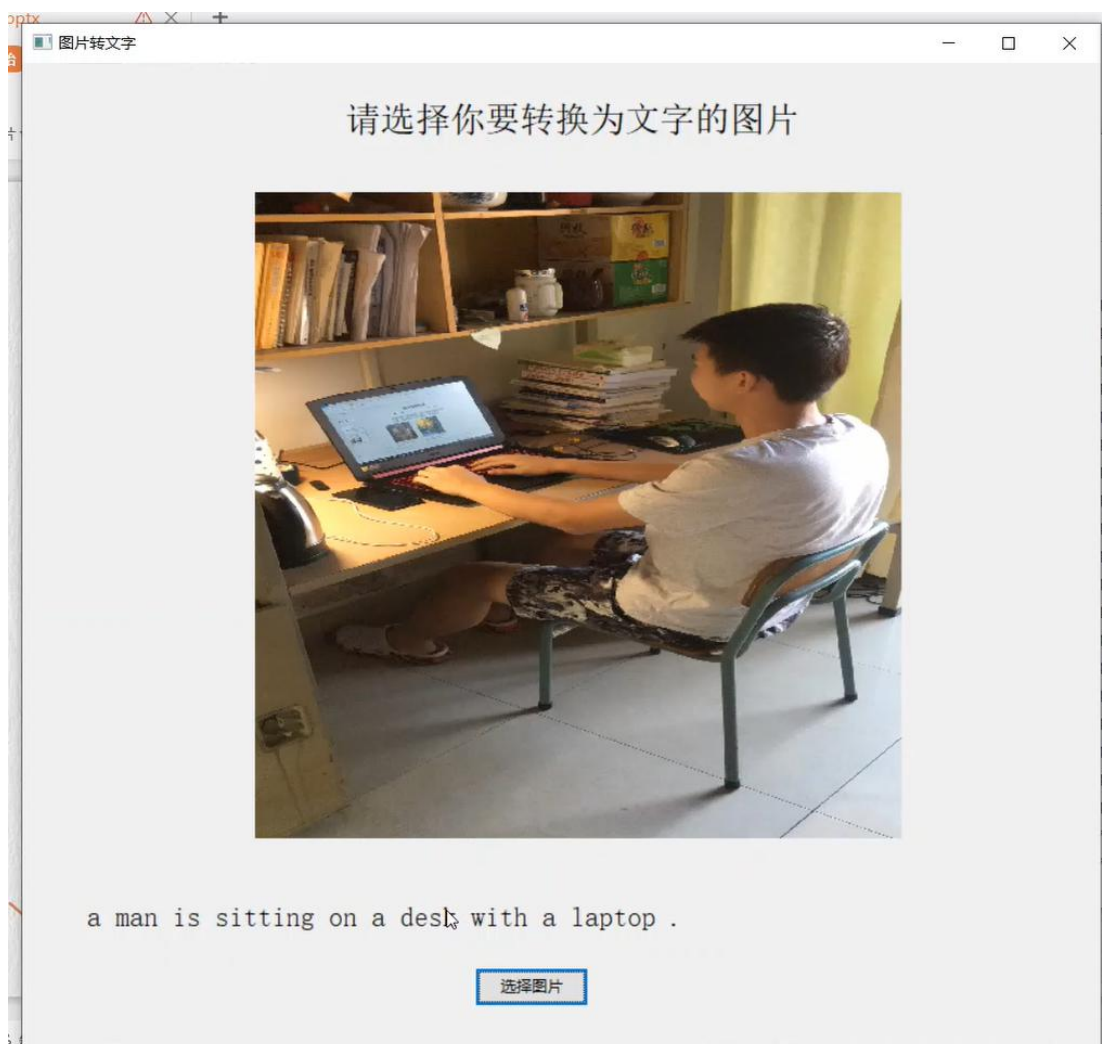
In this paper, we describe approaches to caption generation that attempt to incorporate a form of attention with

相关论文的评价情况

Dataset	Model	BLEU			
		BLEU-1	BLEU-2	BLEU-3	BLEU-4
COCO	CMU/MS Research (Chen & Zitnick, 2014) ^a	—	—	—	—
	MS Research (Fang et al., 2014) ^{†a}	—	—	—	—
	BRNN (Karpathy & Li, 2014) ^o	64.2	45.1	30.4	20.3
	Google NIC ^{†oΣ}	66.6	46.1	32.9	24.6
	Log Bilinear ^o	70.8	48.9	34.4	24.3
	Soft-Attention	70.7	49.2	34.4	24.3
	Hard-Attention	71.8	50.4	35.7	25.0

100% | 633/633 [04:59<00:00, 2.48it/s] BLEU SCORE: 0.6840128708252958, 0.44712500625246043, 0.28219346445189475, 0.18541787110010358]
100% | 633/633 [04:59<00:00, 2.11it/s]

下面是一些图片描述的结果:



请选择你要转换为文字的图片



a baseball player is swinging a bat at a baseball game .

选择图片

请选择你要转换为文字的图片



a woman is holding a cell phone in her hand .

选择图片

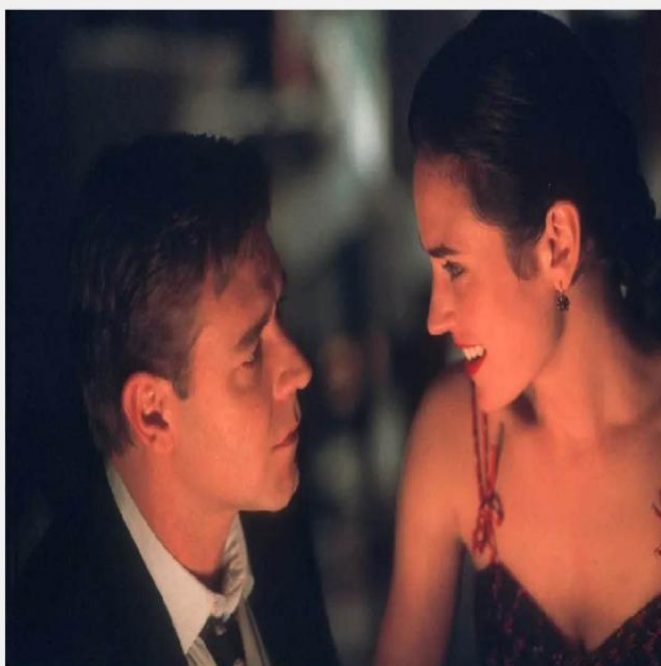
请选择你要转换为文字的图片



a cat sitting on a bench in a room .

选择图片

请选择你要转换为文字的图片



a man and a woman in a suit and tie .

选择图片

请选择你要转换为文字的图片



a group of people sitting on a bench with a woman .

选择图片

总体而言，软件界面简单，功能专一（即图像描述）。图片描述得不算特别精确，但大体描述都还过得去，部分描述比较差劲，当然也有部分描述得相当好，这受限于训练模型的数据集以及训练次数等，毕竟采用的数据集 MSCOCO 里面的事物以及对应词汇量有限；另外，由于小组电脑性能不太好，训练时间太长，所以训练次数相应少很多。但小组算是完成了预期目标，并且有点成果。在这个过程中小组成员付出了许多，也收获了许多，比如接触了深度学习相关知识以及体会了团委项目的整个过程，小组从开始的对项目不知所措到现在总算有点头绪。十分感谢学校及学院提供的平台以及资源，感谢指导老师的珍贵建议与教导，小组成员定会铭记在心，继续前行！