



北京大学

硕士研究生学位论文

题目： 基于像素级处理技术的视频目
标跟踪算法研究

姓 名： 崔家梁
学 号： 1601210300
院 系： 地球与空间科学学院
专 业： 摄影测量与遥感
研究方向： 数字摄影测量与遥感数字成像
导 师： 赵红颖 副教授

2019 年 5 月

版权声明

任何收存和保管本论文各种版本的单位和个人，未经本论文作者同意，不得将本论文转借他人，亦不得随意复制、抄录、拍照或以其他方式传播。否则一旦引起有碍作者著作权之问题，将可能承担法律责任。

摘要

这是中文摘要。
这是中文摘要的第二段。

关键词：数字视频处理，目标跟踪，像素级别

A Research of Video Object Tracking Altorithm based on Pixel-wise Processing Technology

Cui Jialiang (Photogrammetry and Remote Sensing)

Directed by Prof. Zhao Hongying

ABSTRACT

Test of the English abstract.
second parg.

KEYWORDS: Digital video processing, Object tracking, Pixel-level

目录

第一章 引言	1
1.1 研究背景	1
1.1.1 视频目标跟踪问题	1
1.1.2 像素级图像与视频处理技术	1
1.2 研究意义	2
1.2.1 选题来源	2
1.2.2 研究视频跟踪的重要性	2
1.2.3 现阶段视频跟踪算法缺陷	3
第二章 研究现状	4
2.1 视频跟踪算法	4
2.2 像素级别处理算法	5
2.3 像素级别视频目标跟踪算法	5
2.4 深度学习技术	5
2.5 相关数据集	6
第三章 基于 CNN 和 RNN 的像素级别跟踪算法	7
3.1 算法结构与核心思想	7
3.2 基于 CNN 和 RNN 的像素级别跟踪模型在空间维度的处理	7
3.2.1 CNN 的原理	7
3.2.2 多尺度思想的引入	8
3.3 基于 CNN 和 RNN 的像素级别跟踪模型在时间维度的处理	8
3.3.1 RNN 的原理	8
3.3.2 跟踪状态	8
3.3.3 跟踪系统初始化	8
3.4 模型训练	8
3.4.1 机器学习思想	8
3.4.2 模型训练原理	8
第四章 基于 Tensorflow 的像素级视频目标跟踪实验	9
4.1 实验总体设计思路	9
4.2 实验软硬件环境	9

4.2.1 软件环境	9
4.2.2 硬件环境	9
4.3 实验数据	9
4.4 实验程序的编写	10
4.5 实验结果的评估方式	10
第五章 结论与讨论	11
5.1 实验结论	11
5.1.1 实验结果图像展示	11
5.1.2 实验结果定量评估	11
5.2 总结与讨论	11
5.2.1 本研究的创新点	11
5.2.2 本研究的不足	11
5.2.3 后续工作	11
参考文献	12
致谢	14
北京大学学位论文原创性声明和使用授权说明	15

第一章 引言

本章将介绍本研究^①的研究背景与研究意义。

1.1 研究背景

本节将介绍视频目标跟踪，像素级处理，深度学习 (TODO) 等基础概念。这将为后文详细介绍本研究打下基础。

1.1.1 视频目标跟踪问题

视频 (Video; 影片 [港澳台]) 可以看作一个图像序列^②，由连续曝光获取的图像组成。很多情况下，尤其是利用摄像机在对特定的对象进行观察时，视频中将会长期存在一个或多个需要关注的目标。典型的会出现在视频中的目标可以是人物，车辆，动物等。如何使用计算机从视频中高效地持续锁定这些目标是计算机视觉、遥感等研究者十分关注的问题。

在无人机应用场景中，视频目标跟踪是一项很重要的应用。广义上的无人机目标跟踪包括无人机对目标进行识别，定位，追踪的整套过程。本研究关注的视频目标跟踪主要指这其中利用计算机对视频中的目标进行持续的定位这一过程。

1.1.2 像素级图像与视频处理技术

在遥感应用中，对地物的识别需要用到各种各样的图像和视频处理技术。先以图像处理为例，一种图像处理技术的输入通常都是图像，而输出可以有很多种。在算法层面，图像处理技术得到的结果通常是标签 + 位置。如在遥感地物分类的任务中，算法给出的结果是每个像元的类型归属；在无人机目标跟踪任务中，算法给出的结果是目标所在的大致位置。根据输出形式的不同，常见的图像和视频处理算法的结果主要有像素级，矩形级和图像级等表现形式。

像素级 (Pixel-wise) 的处理算法输出结果是一副和原图像尺寸大小几乎相同的图像^③；矩形级 (Box-wise) 的图像处理算法的输出结果通常是一个目标的外包矩形 (Bounding Box) 和目标所属的标签，不需要精确到每个像素，只关注兴趣目标的大致位置和大致形状即可；图像级的图像处理算法指给整张图片贴上一些标签，如图像中有没有某个目标，有哪些类型的目标等等。

① 指基于像素级处理技术的视频目标跟踪算法研究，后文将继续称之为本研究

② 本研究中不考虑视频中的音频信息

③ 在某些特殊情况下，如滤波算法，会丢失一些图像边缘信息，得到的结果可能比原图小一些

像素级的图像处理算法的典型应用有图像分割 (Image Segmentation), 遥感地物分类等. 这些应用需要关注到图像的每一个像素, 使每一个像素有唯一的明确的归属, 并且这些应用十分关心图像中不同类型的区域的边界线. 矩形级的图像处理算法的典型应用有图像目标检测 (Image Object Detection) 和矩形框视频目标跟踪等. 由于最终的需求只关心目标的大致形状和位置, 因此在算法的设计上通常会放弃一些像素级的特征. 图像级的图像处理算法则处理更大尺度的问题. 通常只需宏观上得到正确的结果即可.

1.1.2.1 像素级视频目标跟踪问题

目前主流的视频目标跟踪研究都是在试图解决矩形框视频目标跟踪问题, 且大多数算法设计之初就是为产生矩形级结果而设计, 其中几乎没有提取像素级结果的过程. 也有少数像素级跟踪算法, 但没有用到最新的图像处理技术.

区别于现在大多数矩形级的目标跟踪算法, 本文提出的像素级目标跟踪算法将直接以得到精确的像素级的跟踪结果为目标.

1.2 研究意义

1.2.1 选题来源

在我做本科毕业设计^[1]时, 就曾考虑到无人机影像处理过程中缺乏能将兴趣目标精确提取出的方法. 随着研究的深入, 我发现近年来在计算机视觉的研究中常用的深度学习方法十分适合解决这个问题. 然而将计算机视觉中的视频跟踪的方法主要针对自动驾驶等领域, 将其直接运用于无人机航拍的视频效果并不好, 并且大多数方法无法达到遥感所需要的像素级别处理的需求. 因此我选择了研究针对视频中目标的像素级别跟踪这一问题. 在工业界推动下计算机学者们已经研究出很多成熟的像素级图像处理算法, 将其目标稍加改动, 即可运用在视频, 从而有希望在无人机影像处理中得到较好的效果.

1.2.2 研究视频跟踪的重要性

视频是无人机遥感应用的重要信息载体. 在需要持续观察并关注的遥感应用场景中, 视频由于其较容易获取, TODO go on

由于现阶段的无人机平台已经实现轻量化, 用无人机来跟踪目标是理所当然的最佳选择. 然而现在用于无人机的跟踪算法依然不够强劲. 现有的跟踪算法大多只能在高功率的 PC 上运行, 并且想获得好的跟踪效果就要加大模型, 增加功耗. 因此还无法向无人机平台迁移.

1.2.3 现阶段视频跟踪算法缺陷

现有的跟踪算法为了达到跟踪效果, 通常需要结合多种跟踪方法, 并进行结果的融合. 从算法层面看, 这样的结构既不高校, 也不美观. 实际上, 由于根本理论方法的缺乏, 一些算法不得不设计得越来越复杂. 这种情况下需要新的思路来打破局面.

视频目标跟踪算法由于要面临视频时间和空间纬度的大量数据, 单位时间接受到的信息量极大. 目前的多种跟踪方式均无法准确的从这些信息中提取到最少量的有效信息. 算法的质的提升任然需要理论的创新. 本文研究的主要是像素级别的目标跟踪. 通过研究像素级别的目标跟踪, 或许能获得对现有的外包矩形目标跟踪算法理论上的帮助, 让产生式跟踪模型 (第 2 章中将会介绍) 重新受到重视.

第二章 研究现状

本章将介绍像素级视频影像跟踪算法及其相关算法的研究现状, 为下一章介绍本研究的理论创新铺垫基础。

需要提前指出的是, 本文所提出方法将是像素级别的处理技术在视频目标跟踪问题上的一个应用, 并不是现在狭义上定义的视频目标跟踪算法. 但其中很多思想借鉴了现在的视频目标跟踪算法. 因此在研究现状部分我们依然会着重分析视频目标跟踪算法, 同时将介绍像素级别处理技术与深度学习思想.

2.1 视频跟踪算法

这里介绍的视频目标跟踪算法均是矩形级算法. 视频目标跟踪是目前视频处理的一个很热门的研究方向. 受限于目前计算机的计算能力, 我们不能随意增大算法的规模, 因为在大多数情况下不能实时进行目标跟踪的方法是没有意义的. 因此视频目标跟踪算法必须要节约计算资源. 因而算法的设计就显得格外重要.

视频跟踪算法主要分为产生式模型和判别式模型.

产生式模型指基于当前时刻及前一段时间的目标状态, 结合新加入的帧的视频内容, 直接根据概率模型产生一个新的跟踪目标. 在计算能力极差的八九十年代, 许多早期的模型^[2] 都是产生式模型. 直到 20 世纪初, 产生式模型依然是主流. 基于 Kalman 滤波的许多模型^[3-5] 都为推动跟踪效果做出过贡献.

然而在现在 (2018 年), 判别式模型已经完全占据了视频目标跟踪的主流. 2012 年 Hinton 提出 AlexNet^[6] 后, 深度学习这一划时代的思想迅速站上了图像处理界的主流. 由于卷积神经网络^[6] (Convolution neural network) 在图像处理的普适性, 在图像分类^[6-8], 图像分割^[9] 和目标检测^[10,11] 等方面均赢得了学界的认可, 迅速与传统方法结合, 成为这些研究方向必不可少的重要方法. 在视频跟踪问题上, 深度学习方法同样有较好的表现. 经过几年发展, 脱颖而出的基于深度学习的视频目标跟踪算法主要都是判别式模型. 判别式模型指分两步完成跟踪的一种模型, 第一步是利用提取特征的方法, 将新帧作为一个图像做特征提取运算; 第二步是结合提取出的特征和之前的跟踪结果, 在提取出的特征中选择要跟踪的目标. 具有代表性的有 2016 年的 MDNet^[12] 算法. 该算法的主要思想是利用一个预先训练好的深度神经网络将送入的新帧作为图像提取特征, 再形成多个次级网络进行目标跟踪. 还有一些基于检测的目标跟踪, 如 ROLO^[13] 算法, 先利用目标检测技术检测出很多目标, 再从这些目标中选择一个和正在跟踪的目标比较像的目标作为跟踪结果.

如前文所说, 由于当今计算能力的爆发, 由于能很快得到大量的目标检测结果, 判别式模型大行其道. 但判别式模型从思想上是目标检测的产物, 其执行过程中将花费大量的时间去生成根本不是被跟踪目标的其他目标. 本文将提出的是一种生成式模型, 试图重新从目标跟踪的本质任务出发.

2.2 像素级别处理算法

区别于典型的图像分类与目标外包框检测问题, 像素级别 (Pixel-wised) 的图像处理需要获得一个覆盖全图的, 精确到目标轮廓信息的结果. 现在最常见的像素级别应用是图像分割. 在图像分割领域, 以早期的分水岭算法^[14]为代表的传统阵营^①已经有一系列研究. 虽然分水岭算法有许多改进^[15], 但只能在大尺度图像上表现较好. 对复杂情况下的分割效果依然不够智能. 在深度学习技术出现后, 深度学习很快就被运用于分割领域. 最成功的典型是由 U-Net^[16] 开创的降级-升级模型. 与之类似的还有 SegNet^[17] 将 U-Net 的升级模型稍加改动后得到了更好的效果.

但这些算法都是为图像分割设计的. 相比与图像处理算法, 视频处理算法更需要注意帧与帧之间的关系. 上一段提到的著名的 SegNet 图像分割算法的演示阶段用的是视频数据做展示, 但其只是将视频拆成了完全独立的图像进行处理, 仔细观察会发现许多细节的处理会缺乏连贯性.

2.3 像素级别视频目标跟踪算法

近两年来, 像素级别的视频目标跟踪算法也有所研究. 如纽约大学在 2017 年完成的一项工作^[18], 该文用 Conv-LSTM 技术^[19] 尝试了对像素级别目标的视频目标跟踪. 但该研究使用的手段复杂, 最终结合了两种跟踪算法才得到结果. 在更早的 2007 年, Hua 等人用 K-means 算法也尝试过像素级别的跟踪^[20], 但由于没有结合深度学习算法, 得到的结果也并不理想.

本文希望提出的像素级别跟踪算法将建立在一个单独简洁的框架上, 结合目前像素级别处理技术和目标跟踪技术的精髓, 实现一个思路清晰的算法, 并尝试寻求更好的结果.

2.4 深度学习技术

2012 年 Hinton 等人主导的深度学习技术能在图像处理, 自然语言处理等方向大放异彩的主要原因之一, 是深度学习采用了简单的模型, 配合上复杂而可训练的参数, 从

① 这里的传统阵营指用非神经网络方法的算法的处理方式

而得到更好的结果. 简单的模型指 CNN^②, RNN^③. 前者主要用于处理空间尺度, 即图像; 后者主要处理时间尺度, 即语音和视频. CNN 和 RNN 的结构都很简单, 由于权值共享思想, 需要训练的参数也很少. 但由于其采用了仿生学的原理, 得到的结果往往比传统算法优秀. 实际上在 20 世纪末, 就已经有人提出并用图像数据尝试了深度学习^[21]. 但直到近年来, GPU 计算的普及才使得深度学习技术有了更大的用武之地. 2013 年, 加州大学伯克利分校的发布了 Caffe^[22] 深度学习工具 (后来与 PyTorch^[23] 合并), 谷歌公司于 2017 年初发布 Tensorflow^[24] 的 python 版本 API. 这些开源, 开放, 高效, 简单易用的工具使深度学习算法的实现变得十分容易.

本文提出的算法将主要采用深度学习方法, 力求用一个简单的深度学习模型解决复杂的问题.

2.5 相关数据集

近年来随着研究的火热, 产生了许多网络上共享的数据集, 典型的有 2009 年的 ImageNet^[25]. 由于深度学习需要大量的训练数据, 开放的数据集直接推动了深度学习的发展.

在跟踪领域最典型的有 VOT^[26] 和 OTB^[27] 数据集. 特别是 VOT 数据集 2016 年的像素级别数据^[28], 以人工标注的方式提供了像素级别的视频跟踪训练集. 在图像分割问题上同样有许多数据集, 如 VOS^[29] 等. 这些数据集的数据量很大, 数据质量也很好, 给模型训练带来了许多方便. 本文将直接使用 VOT 等跟踪数据集, 并尝试使用一些分割数据集对模型进行更为细节的训练.

② Convolution Neural Network, 卷积神经网络

③ Recurrent Neural Network, 循环神经网络

第三章 基于 CNN 和 RNN 的像素级别跟踪算法

本章是本文的重点, 将详细介绍本研究的理论创新。

本研究提出一种结合现有的像素级别处理技术和现有的矩形级视频目标跟踪技术的像素级别目标跟踪算法, 具体算法结构, 实现细节, 训练等将在本章重点介绍。

3.1 算法结构与核心思想

本文所实现的算法将首先基于用于实现静态图像分割的 U-Net^[16] 的多级降级-升级卷积神经网络结构, 但将在这个多级网络结构中加入 Conv-LSTM 结构。

类似与 U-Net 的结构, 本文的卷积网络部分也将有多个降级和升级结构; 每个降级结构包括几个卷积层, 使用池化结构进行降级; 每个升级部分采用升卷积进行升级处理。在降级过程中, 图片数据的尺寸大小会衰减, 同时等比例增加其波段范围。对于 3 层的结构, 最小级的波段将有 128 个。这个多级结构的设计理念是为了处理多尺度问题; 浅层的级别能很好的处理细节问题, 但对宏观的把控会较弱, 具体表现为可能会出现噪声点; 深层的结构对宏观把控好, 但对边界处理较弱。升级结构能将浅层处理得到的边界信息与深层处理得到的宏观信息相结合, 得到一个更好的结果。

在时间尺度, 本研究的算法将主要采用 LSTM 算法解决问题。具体的, LSTM 单元将被加入到各个层级当中。LSTM 在各种跟踪算法中有广泛应用, 但大多数算法仅仅将其作为对最后结果的处理手段。本研究的算法将把 LSTM 作为所有的中间状态记录单元。

与纽约大学 2017 年实现的 Conv-LSTM 结构的跟踪算法不同的是, 本文所采用的多级神经网络将把 Conv-LSTM 加入各个卷积层级; 而与 U-Net, SegNet 等多级分割算法不同的是, 本文将在整个结构中多处穿插 LSTM 以得到一个时间连续的结果。

3.2 基于 CNN 和 RNN 的像素级别跟踪模型在空间维度的处理

3.2.1 CNN 的原理

CNN 是一种基于模板滤波的图像处理方法。常用的 CNN 的模板大小是 3×3 , 有时也可以是 5×5 , 7×7 等大小

3.2.2 多尺度思想的引入

3.3 基于 CNN 和 RNN 的像素级别跟踪模型在时间维度的处理

3.3.1 RNN 的原理

3.3.2 跟踪状态

3.3.3 跟踪系统初始化

3.4 模型训练

3.4.1 机器学习思想

3.4.2 模型训练原理

第四章 基于 Tensorflow 的像素级视频目标跟踪实验

为了验证本研究提出的方法的可行性和价值,我们设计了一个实验。后文中该实验将被称为本实验。本实验基本实现了本研究提出的方法,并得到了一定的结果和结论。本章将介绍本实验的实验条件与过程。

本章接下来的部分将介绍本实验的设计思路,硬件环境和数据等实验条件,实验代码的实现和实验结果的评估方法。

4.1 实验总体设计思路

4.2 实验软硬件环境

4.2.1 软件环境

本实验的软件部分主要在 Tensorflow^[24] 框架下实现。

Tensorflow 是最初由谷歌公司开发的一套现以开源的机器学习框架,可以为算法研究者屏蔽操作系统与硬件,资源分配,梯度计算等繁琐部分,让研究者能将更多的注意力集中在算法过程中。对于本研究,Tensorflow 主要贡献了 CNN,RNN 单元的结构定义,损失函数定义,正向反向传播与梯度更新等功能。

4.2.2 硬件环境

本实验几乎所有的运算操作是在一台配置有英伟达 GTX1070 图形处理器,英特尔 i7 中央处理器,24GB 内存的笔记本电脑上进行的。

本实验深度学习计算部分使用了 GPU 加速,直接依赖 Tensorflow 的 GPU 选项进行。本研究曾尝试过只用 CPU 进行计算,也能得到一定结果。

如果有更好的硬件条件(更多、更好的图形处理器,更大的内存,更多核心的 CPU),本实验有希望会得到更精细的结果。

4.3 实验数据

本实验使用 VOT2016 数据集^[28] 实现,相似的数据集还有 VOT2017 等。

如图4.1所示,该数据集通过人工标记,提供了十分优秀的像素级的目标跟踪数据。该数据集有几百个序列,共有几万张图片。

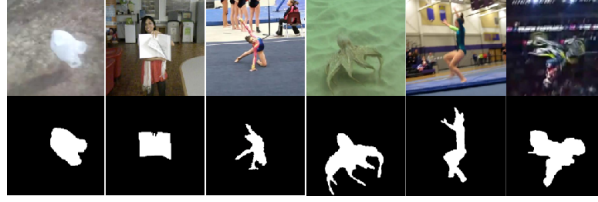


图 4.1 VOT2016 像素级标记

本实验的训练集和测试集均来源于该数据集，使用时将所有数据随机切分为训练集和测试集。

4.4 实验程序的编写

事实上, 虽然借助于 Tensorflow 实现了许多计算功能, 但本研究依然经历了许多代码开发工作, 包括但不限于神经网络结构定义, 训练数据处理等.

4.5 实验结果的评估方式

第五章 结论与讨论

5.1 实验结论

5.1.1 实验结果图像展示

5.1.2 实验结果定量评估

5.2 总结与讨论

5.2.1 本研究的创新点

5.2.2 本研究的不足

5.2.3 后续工作

参考文献

- [1] 崔家梁. 针对无人机复杂畸变视频影像的稳像后处理方法. **2016**.
- [2] Robert J Schalkoff and Eugene S Mcvey. “A model and tracking algorithm for a class of video targets”. *IEEE transactions on pattern analysis and machine intelligence*, **1982**(1): 2–10.
- [3] Changick Kim and Jenq-Neng Hwang. “Fast and automatic video object segmentation and tracking for content-based applications”. *IEEE transactions on circuits and systems for video technology*, **2002**, 12(2): 122–129.
- [4] Shiuh-Ku Weng, Chung-Ming Kuo and Shu-Kang Tu. “Video object tracking using adaptive Kalman filter”. *Journal of Visual Communication and Image Representation*, **2006**, 17(6): 1190–1208.
- [5] Dorin Comaniciu, Visvanathan Ramesh and Peter Meer. “Kernel-based object tracking”. *IEEE Transactions on pattern analysis and machine intelligence*, **2003**, 25(5): 564–577.
- [6] Alex Krizhevsky, Ilya Sutskever and Geoffrey E Hinton. “Imagenet classification with deep convolutional neural networks”. In: *Advances in neural information processing systems*. **2012**: 1097–1105.
- [7] Ian H Witten, Eibe Frank, Mark A Hall *et al.* *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, **2016**.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren *et al.* “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. **2016**: 770–778.
- [9] Jonathan Long, Evan Shelhamer and Trevor Darrell. “Fully convolutional networks for semantic segmentation”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. **2015**: 3431–3440.
- [10] Shaoqing Ren, Kaiming He, Ross Girshick *et al.* “Faster r-cnn: Towards real-time object detection with region proposal networks”. In: *Advances in neural information processing systems*. **2015**: 91–99.
- [11] Joseph Redmon, Santosh Divvala, Ross Girshick *et al.* “You only look once: Unified, real-time object detection”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. **2016**: 779–788.
- [12] Hyeonseob Nam and Bohyung Han. “Learning Multi-Domain Convolutional Neural Networks for Visual Tracking”. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016-06.
- [13] Guanghan Ning, Zhi Zhang, Chen Huang *et al.* “Spatially Supervised Recurrent Convolutional Neural Networks for Visual Object Tracking”. *arXiv preprint arXiv:1607.05781*, **2016**.
- [14] Ole Fogh Olsen and Mads Nielsen. “Multi-scale gradient magnitude watershed segmentation”. In: *International Conference on Image Analysis and Processing*. **1997**: 6–13.
- [15] Vicente Grau, AUJ Mewes, M Alcaniz *et al.* “Improved watershed transform for medical image segmentation using prior information”. *IEEE transactions on medical imaging*, **2004**, 23(4): 447–458.

-
- [16] Olaf Ronneberger, Philipp Fischer and Thomas Brox. “U-net: Convolutional networks for biomedical image segmentation”. In: *International Conference on Medical image computing and computer-assisted intervention*. **2015**: 234–241.
- [17] Vijay Badrinarayanan, Alex Kendall and Roberto Cipolla. “Segnet: A deep convolutional encoder-decoder architecture for image segmentation”. *IEEE transactions on pattern analysis and machine intelligence*, **2017**, 39(12): 2481–2495.
- [18] Yilin Song, Chenge Li and Yao Wang. “Pixel-wise object tracking”. *CoRR*, **2017**, abs/1711.07377. <http://arxiv.org/abs/1711.07377>.
- [19] Viorica Ptrucean, Ankur Handa and Roberto Cipolla. “Spatio-temporal video autoencoder with differentiable memory”. In: *International Conference on Learning Representations (ICLR) Workshop*. **2016**.
- [20] Chunsheng Hua, Haiyuan Wu, Qian Chen *et al.* “K-means clustering based pixel-wise object tracking”. *Information and Media Technologies*, **2008**, 3(4): 820–833.
- [21] Yann LeCun, Léon Bottou, Yoshua Bengio *et al.* “Gradient-based learning applied to document recognition”. *Proceedings of the IEEE*, **1998**, 86(11): 2278–2324.
- [22] Yangqing Jia, Evan Shelhamer, Jeff Donahue *et al.* “Caffe: Convolutional Architecture for Fast Feature Embedding”. *arXiv preprint arXiv:1408.5093*, **2014**.
- [23] Adam Paszke, Sam Gross, Soumith Chintala *et al.* “Automatic differentiation in PyTorch”. **2017**.
- [24] Martin Abadi, Paul Barham, Jianmin Chen *et al.* “TensorFlow: A System for Large-Scale Machine Learning.” In: *OSDI*. **2016**: 265–283.
- [25] J. Deng, W. Dong, R. Socher *et al.* “ImageNet: A Large-Scale Hierarchical Image Database”. In: *CVPR09*. **2009**.
- [26] Matej Kristan, Jiri Matas, Ale Leonardis *et al.* “A Novel Performance Evaluation Methodology for Single-Target Trackers”. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016-11, 38(11): 2137–2155.
- [27] Yi Wu, Jongwoo Lim and Ming-Hsuan Yang. “Online Object Tracking: A Benchmark”. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. **2013**.
- [28] Tomas Vojir and Jiri Matas. *Pixel-Wise Object Segmentations for the VOT 2016 Dataset* [Research Report]. Prague, Czech Republic, 2017-01.
- [29] S. Caelles, K.K. Maninis, J. Pont-Tuset *et al.* “One-Shot Video Object Segmentation”. In: *Computer Vision and Pattern Recognition (CVPR)*. **2017**.

致谢

这是致谢的第一段

这是致谢的第二段

北京大学学位论文原创性声明和使用授权说明

原创性声明

本人郑重声明：所呈交的学位论文，是本人在导师的指导下，独立进行研究工作所取得的成果。除文中已经注明引用的内容外，本论文不含任何其他个人或集体已经发表或撰写过的作品或成果。对本文的研究做出重要贡献的个人和集体，均已在文中以明确方式标明。本声明的法律结果由本人承担。

论文作者签名： 日期： 年 月 日

学位论文使用授权说明

（必须装订在提交学校图书馆的印刷本）

本人完全了解北京大学关于收集、保存、使用学位论文的规定，即：

- 按照学校要求提交学位论文的印刷本和电子版本；
- 学校有权保存学位论文的印刷本和电子版，并提供目录检索与阅览服务，在校园网上提供服务；
- 学校可以采用影印、缩印、数字化或其它复制手段保存论文；
- 因某种特殊原因需要延迟发布学位论文电子版，授权学校在 ☐ 一年 / ☐ 两年 / ☐ 三年以后在校园网上全文发布。

（保密论文在解密后遵守此规定）

论文作者签名： 导师签名： 日期： 年 月 日