

HOW TO WRITE THE INTRODUCTION FOR A RESEARCH PAPER

in three easy steps

Kate Saenko

INTRODUCTION

This talk will show how to write an introduction to a research paper in 3 easy steps

- STEP1: make a plan
- STEP2: write!
- STEP3: edit edit edit



WWW.PHDCOMICS.COM

STEP 1: MAKE A PLAN

STEP 1: MAKE A PLAN

Write down **in bullets**:

- your contribution
- what problem is it solving?
- what is the evaluation that proves it?

STEP 1: MAKE A PLAN

EXAMPLE

Contribution:

- Propose dense caption task
- Design new end-to-end model

Problem:

- Current methods fail...
- No one has tried this...

Evaluation:

- Compare to baseline X on dataset Y
- Tease a promising result

Kate Saenko

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in images in natural language. The dense captioning task generalizes object detection when the descriptions consist of a single word, and Image Captioning when one predicted region covers the full image. To address the localization and description task jointly we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Our ability to effortlessly point out and describe all aspects of an image relies on a strong semantic understanding of a visual scene and all of its elements. However, despite numerous potential applications, this ability remains a challenge for our state of the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to assign one label to an image. Further work has pushed these advances along two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified models that efficiently identify and label multiple salient regions of an image. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequence of words able to express significantly richer concepts.

However, despite encouraging progress along the label density and label complexity axes, these two directions have

* Both authors contributed equally to this work.



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (see Figure 1), which requires a model to predict a set of descriptions across regions of an image. Object detection is hence recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image.

Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region.

We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region captions. Our results show both performance and speed improvements over approaches based on previous state of the art. We make our code and data publicly available to support further progress on the dense captioning task.

ONE IDEA IS BETTER THAN MANY!



Good: reader has an “aha!”
moment



Bad: too many things at once

STEP 1: EXPAND PLAN

Re-write bullets into following structure*:

- task
- state-of-the-art
- flaw in state-of-the-art
- your idea/solution
- proof it works

*this is usually easier to do for beginners than starting with this structure

STEP 1: EXPAND PLAN

EXAMPLE

title

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in an image in natural language. The dense captioning task is more challenging than traditional image captioning when the descriptions are more detailed and localized. To address the localization, we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external region proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

task

SotA

flaw

essly point out and describe all aspects of an image with a strong semantic understanding of all of its elements. However, despite numerous publications, this ability remains a challenge in the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to identify a single class for the whole image. Further work has pushed the state of the art in two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified modern methods and label multiple salient regions. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequences of significantly richer concepts. Encouraging progress along the localization and captioning axes, these two directions have remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (Figure 1), which requires a model to predict a set of localized descriptions across regions of an image. Object detection is recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image. Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region. We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region-grounded captions. Our results show both performance and speed improvements over approaches based on previous art. We make our code and data publicly available to support further progress on the dense captioning task.

figure 1



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

idea

proof

STEP 1: EXPAND PLAN

TASK

What can it do?

- State the task your model solves
- Usually only need one sentence
- Can mention applications

Why do we need this?

- Avoid reader frustration!
- (reads to page 7): oh, they're doing detection, not classification!? \$#@#%\$!!

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in images in natural language. The dense captioning task generalizes object detection when the descriptions consist of a single word, and Image Captioning when one predicted region covers the full image. To address the localization and description task jointly we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external regions proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Essentially point out and describe all aspects of a strong semantic understanding of a task. However, despite applications, this ability remains a challenge of the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to assign one label to an image. Further work has pushed these advances along two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified models that efficiently identify and label multiple salient regions of an image. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequence of words able to express significantly richer concepts.

However, despite encouraging progress along the label density and label complexity axes, these two directions have

* Both authors contributed equally to this work.



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (see Figure 1), which requires a model to predict a set of descriptions across regions of an image. Object detection is hence recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image.

Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region.

We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region captions. Our results show both performance and speed improvements over approaches based on previous state of the art. We make our code and data publicly available to support further progress on the dense captioning task.

STEP 1: EXPAND PLAN

STATE-OF-THE-ART

What is the best we can do now?

- Briefly cite recent work (direct competition)
- Get the reader excited about area
- "Significant progress", "exciting results"

Common mistakes:

- Writing more than one paragraph (that is what related work is for!)
- Diminishing prior work. You're standing on the shoulders of giants. Praise is better.

Kate Saenko

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jjohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in images in natural language. The dense captioning task generalizes object detection when the descriptions consist of a single word, and Image Captioning when one predicted region covers the full image. To address the localization and description task jointly we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external regions proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Our ability to effortlessly point out and describe all aspects of an image relies on a strong semantic understanding of a visual scene and all of its elements. However, despite numerous potential applications, this ability remains a challenge for our state of the art visual recognition systems. In the last few years there has been significant progress in image localization [39, 26, 53, 45], where the task is to find a bounding box that covers the object of interest in an image. Further work has pushed this task in two orthogonal directions: First, rapid detection [40, 14, 46] has identified modernity and label multiple salient regions [4, 3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequence of words able to express significantly richer concepts.

However, despite encouraging progress along the label density and label complexity axes, these two directions have

*Both authors contributed equally to this work.



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (see Figure 1), which requires a model to predict a set of descriptions across regions of an image. Object detection is hence recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image.

Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region.

We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region captions. Our results show both performance and speed improvements over approaches based on previous state of the art. We make our code and data publicly available to support further progress on the dense captioning task.

SotA

STEP 1: EXPAND PLAN

FLAW

What can't we do yet?

- What is the critical flaw in SotA? (you'll be providing a solution to this later...)
- Why should the reader care? No, really?
- Why can't X, Y or Z be used to solve this?

Engage the reader

- Use an example: text, or visual (see Figure 1)
- At the end, the reader should be thinking, wow, this is really a important problem! I really wish someone would try to solve it!!

Kate Saenko

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in images in natural language. The dense captioning task generalizes object detection when the descriptions consist of a single word, and Image Captioning when one predicted region covers the full image. To address the localization and description task jointly we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external regions proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Our ability to effortlessly point out and describe all aspects of an image relies on a strong semantic understanding of a visual scene and all of its elements. However, despite numerous potential applications, this ability remains a challenge for our state of the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to assign one label to an image. Further work has pushed these advances along two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified models that efficiently identify and label multiple salient regions of an image. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequences of significantly richer concepts. Encouraging progress along the label complexity axes, these two directions have remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (see Figure 1), which requires a model to predict a set of descriptions across regions of an image. Object detection is hence recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image. Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region. We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region captions. Our results show both performance and speed improvements over approaches based on previous state of the art. We make our code and data publicly available to support further progress on the dense captioning task.



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

flaw

ated equally to this work.

STEP 1: EXPAND PLAN

IDEA

What you propose to fix the flaw

- The “hero”
- Comes to the rescue

Keep it brief

- Give the main intuition
- “We solve this by adding a new DNN layer”
- Leave the gory details to the method section!

Tip: if you have 2 ideas, describe each ‘flaw’ before you introduce each idea

Kate Saenko

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in images in natural language. The dense captioning task generalizes object detection when the descriptions consist of a single word, and Image Captioning when one predicted region covers the full image. To address the localization and description task jointly we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external regions proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Our ability to effortlessly point out and describe all aspects of an image relies on a strong semantic understanding of a visual scene and all of its elements. However, despite numerous potential applications, this ability remains a challenge for our state of the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to assign one label to an image. Further work has pushed these advances along two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified models that efficiently identify and label multiple salient regions of an image. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequence of words able to express significantly richer concepts.

However, despite encouraging progress along the label density and label complexity axes, these two directions have

*Both authors contributed equally to this work.



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

remained separate. In this work we take a step to fixing these two inter-connected tasks into one joint work. First, we introduce the dense captioning (Figure 1), which requires a model to predict a set of regions across regions of an image. Object detection recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image.

Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region.

We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region captions. Our results show both performance and speed improvements over approaches based on previous state of the art. We make our code and data publicly available to support further progress on the dense captioning task.

idea

STEP 1: EXPAND PLAN

FIGURE 1

“Must have” for a vision paper

- Convey the flaw and idea visually
- Like an “ad” for your paper

Many people ONLY look at this

- If it looks interesting, they will also check the main result table
- Caption should be self-explanatory, since they may not read the main text

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in images in natural language. The dense captioning task generalizes object detection when the descriptions consist of a single word, and Image Captioning when one predicted region covers the full image. To address the localization and description task jointly we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external region proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Our ability to effortlessly point out and describe all aspects of an image relies on a strong semantic understanding of a visual scene and all of its elements. However, despite numerous potential applications, this ability remains a challenge for our state of the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to assign one label to an image. Further work has pushed these advances along two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified models that efficiently identify and label multiple salient regions of an image. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequence of words able to express significantly richer concepts.

However, despite encouraging progress along the label density and label complexity axes, these two directions have

* Both authors contributed equally to this work.



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (see Figure 1), which requires a model to predict a set of descriptions across regions of an image. Object detection is hence recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image.

Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region.

We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region captions. Our results show both performance and speed improvements over approaches based on previous state of the art. We make our code and data publicly available to support further progress on the dense captioning task.

figure 1

STEP 1: EXPAND PLAN

FIGURE 1

OK figure :| illustrates the task well

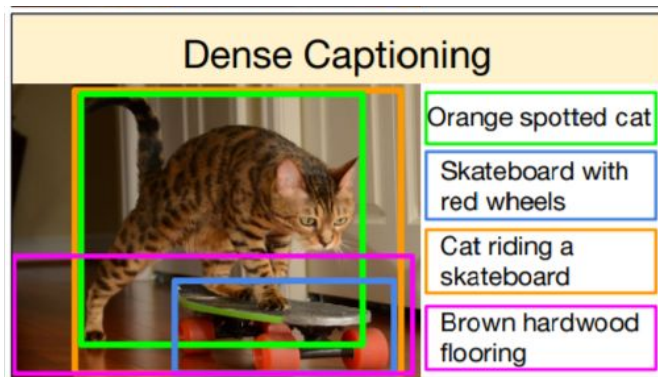


Figure 1. We address the Dense Captioning task with a model that jointly generates both dense and rich annotations in a single forward pass.

Better :) shows flaw in prior work

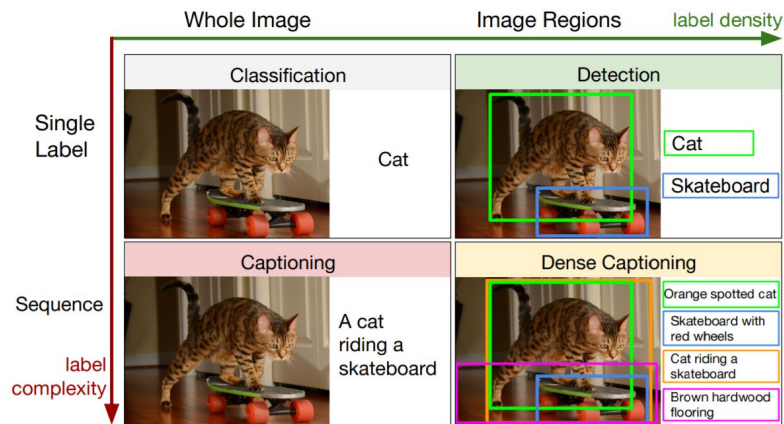


Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

STEP 1: EXPAND PLAN

PROOF

Does it work?

- Say a bit about evaluation and baselines
- Don't give all away, but tease
- E.g. "our method lowers reconstruction error by 20% compared to previous methods"

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in images in natural language. The dense captioning task generalizes object detection when the descriptions consist of a single word, and Image Captioning when one predicted region covers the full image. To address the localization and description task jointly we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external regions proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Our ability to effortlessly point out and describe all aspects of an image relies on a strong semantic understanding of a visual scene and all of its elements. However, despite numerous potential applications, this ability remains a challenge for our state of the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to assign one label to an image. Further work has pushed these advances along two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified models that efficiently identify and label multiple salient regions of an image. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequence of words able to express significantly richer concepts.

However, despite encouraging progress along the label density and label complexity axes, these two directions have

* Both authors contributed equally to this work.



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (see Figure 1), which requires a model to predict a set of descriptions across regions of an image. Object detection is hence recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image.

Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region.

We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region-grounded captions. Our results show both performance and speed improvements over approaches based on previous work. We make our code and data publicly available to support further progress on the dense captioning task.

proof

STEP 1: EXPAND PLAN

TITLE

Specific:

- DenseCap: Fully Convolutional Localization Networks for Dense Captioning :)
- A Network for Describing Images :(
- A Deep Learning Approach Based on Convolutional Neural Network... :(

Memorable:

- If using acronym/abbreviation make sure it is easy to spell and pronounce, e.g. FCLN :(vs DenseCap :)
- Best if conveys central idea, "dense captioning"
- "It's a bird! It's a plane! ..." and other catchphrases: can be corny, proceed with caution !!!

Kate Saenko

title

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in images in natural language. The dense captioning task generalizes object detection when the descriptions consist of a single word, and Image Captioning when one predicted region covers the full image. To address the localization and description task jointly we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external region proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Our ability to effortlessly point out and describe all aspects of an image relies on a strong semantic understanding of a visual scene and all of its elements. However, despite numerous potential applications, this ability remains a challenge for our state of the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to assign one label to an image. Further work has pushed these advances along two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified models that efficiently identify and label multiple salient regions of an image. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequence of words able to express significantly richer concepts.

However, despite encouraging progress along the label density and label complexity axes, these two directions have

* Both authors contributed equally to this work.



Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (see Figure 1), which requires a model to predict a set of descriptions across regions of an image. Object detection is hence recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image.

Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region.

We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region captions. Our results show both performance and speed improvements over approaches based on previous state of the art. We make our code and data publicly available to support further progress on the dense captioning task.

STEP 1: EXPAND PLAN

ABSTRACT

Concise:

- Summary of task and contributions
- Write this last! Copy-paste, rephrase, done!

Audience:

- Meant for general audience
- "Grandmother" should get the gist

DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson* Andrej Karpathy* Li Fei-Fei
Department of Computer Science, Stanford University
{jcojohns, karpathy, feifeili}@cs.stanford.edu

Abstract

We introduce the dense captioning task, which requires a computer vision system to both localize and describe salient regions in an image in natural language. The dense captioning task is a generalization of image captioning when the descriptions are localized to specific regions of the image. To address the localization, we propose a Fully Convolutional Localization Network (FCLN) architecture that processes an image with a single, efficient forward pass, requires no external regions proposals, and can be trained end-to-end with a single round of optimization. The architecture is composed of a Convolutional Network, a novel dense localization layer, and Recurrent Neural Network language model that generates the label sequences. We evaluate our network on the Visual Genome dataset, which comprises 94,000 images and 4,100,000 region-grounded captions. We observe both speed and accuracy improvements over baselines based on current state of the art approaches in both generation and retrieval settings.

1. Introduction

Our ability to effortlessly point out and describe all aspects of an image relies on a strong semantic understanding of a visual scene and all of its elements. However, despite numerous potential applications, this ability remains a challenge for our state of the art visual recognition systems. In the last few years there has been significant progress in image classification [39, 26, 53, 45], where the task is to assign one label to an image. Further work has pushed these advances along two orthogonal directions: First, rapid progress in object detection [40, 14, 46] has identified models that efficiently identify and label multiple salient regions of an image. Second, recent advances in image captioning [3, 32, 21, 49, 51, 8, 4] have expanded the complexity of the label space from a fixed set of categories to sequence of words able to express significantly richer concepts.

However, despite encouraging progress along the label density and label complexity axes, these two directions have

* Both authors contributed equally to this work.

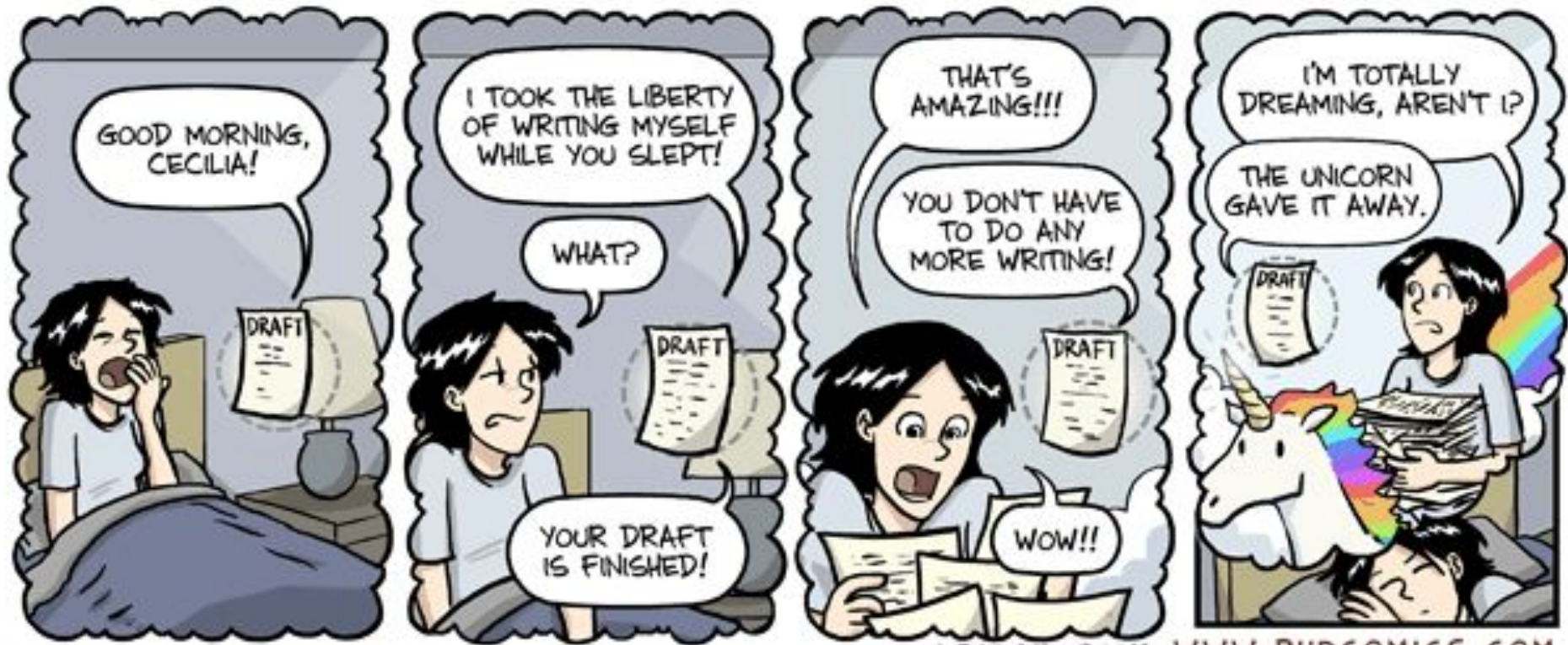


Figure 1. We address the Dense Captioning task (bottom right) with a model that jointly generates both dense and rich annotations in a single forward pass.

remained separate. In this work we take a step towards unifying these two inter-connected tasks into one joint framework. First, we introduce the dense captioning task (see Figure 1), which requires a model to predict a set of descriptions across regions of an image. Object detection is hence recovered as a special case when the target labels consist of one word, and image captioning is recovered when all images consist of one region that spans the full image.

Additionally, we develop a Fully Convolutional Localization Network (FCLN) for the dense captioning task. Our model is inspired by recent work in image captioning [49, 21, 32, 8, 4] in that it is composed of a Convolutional Neural Network and a Recurrent Neural Network language model. However, drawing on work in object detection [38], our second core contribution is to introduce a new dense localization layer. This layer is fully differentiable and can be inserted into any neural network that processes images to enable region-level training and predictions. Internally, the localization layer predicts a set of regions of interest in the image and then uses bilinear interpolation [19, 16] to smoothly crop the activations in each region.

We evaluate the model on the large-scale Visual Genome dataset, which contains 94,000 images and 4,100,000 region captions. Our results show both performance and speed improvements over approaches based on previous state of the art. We make our code and data publicly available to support further progress on the dense captioning task.



JORGE CHAM © 2016 WWW.PHDCOMICS.COM

STEP 2: WRITE!

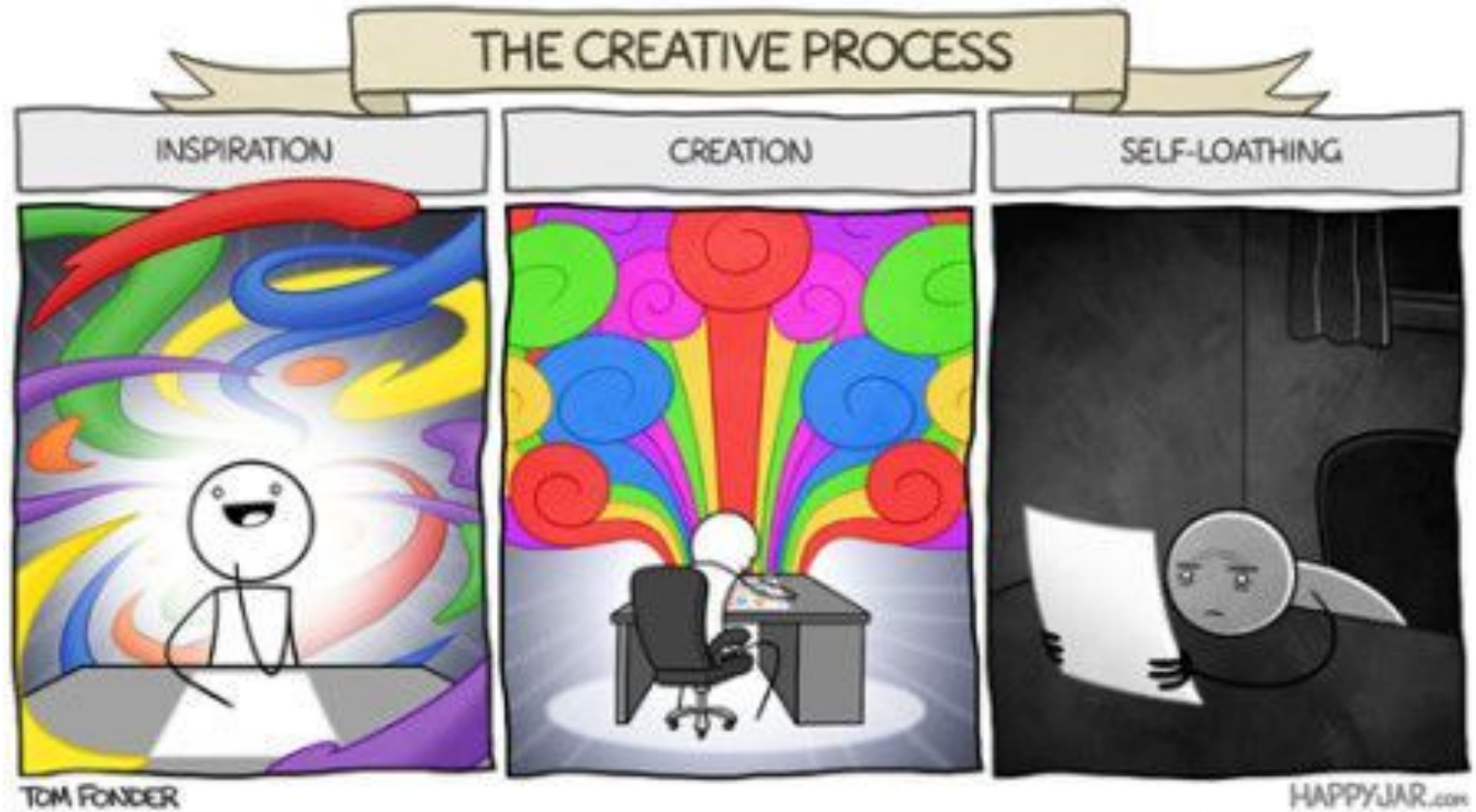
STEP 2: WRITE!

By now, you should have your introduction outlined in several dozen bullets, so the rest should be much easier

This would be a good time to run it by others

Expand each bullet into sentence/paragraph

Fill in details along the way



STEP 3: EDIT EDIT EDIT

STEP 3: EDIT EDIT EDIT

Some lucky people can write near-perfect text in one go

Most of us are not so lucky

I re-write my drafts many times, sometimes dozens

Important to get feedback from your advisor/labmates

THAT'S IT!

Now you just have the rest of the paper to write :)