IMPROVING VISUAL QUESTION ANSWERING USING ACTIVE PERCEPTION ON STATIC IMAGES

Theodoros Bozinis, Nikolaos Passalis and Anastasios Tefas Aristotle University of Thessaloniki (Greece) email: mpozinit@csd.auth.gr, passalis@csd.auth.gr, tefas@csd.auth.gr

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 871449 (OpenDR).

Introduction

- Visual Question Answering (VQA) very challenging deep learning application
- Visual attention key part of VQA

 Correctly identify the region of an image
 Relevant to the question

• Existing methods

Image analysis on fixed & low-resolution images

- \odot Losing fine-grained details
- \circ Sensitive on object scale

Our proposal

- Reinforcement learning-based active perception approach
- Transformation operations on the images (zoom & translation)
- Allows us to
 - perform fine-grained visual analysis
 - effectively increasing the resolution at which the models process information
- Orthogonal to existing attention mechanisms
- Can be combined with existing VQA methods

Virtual camera

Fixed resolution virtual camera

- Keeping the resolution of the image analysis fixed at that
- Not increasing the computational cost
- a) Fine-grained analysis
- b) Keeping only the information relevant to the question
- c) Mitigating the effect of object scale



The setup (1/3)

- input image: $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$
- encoded question: $\boldsymbol{q} \in \mathbb{R}^{N_W}$
- training set $\mathcal{X} = \{ (\mathbf{x}_i, \mathbf{q}_{i,}, \mathbf{t}_i) | \forall i = 1, ..., N \}$
- VQA model: $f_W(\cdot) \in \mathbb{R}^{N_c}$
- $W = argmin_i \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(f_{W'}(\boldsymbol{x}_i, \boldsymbol{q}_i), \boldsymbol{t}_i)$ where $\mathcal{L}(\cdot)$ cross-entropy loss closed set of N_c possible answers

- action space \mathcal{A} a set of 7 actions:
 - a_{left} , a_{right} , a_{up} , a_{down} translation transformations of the camera by δ_{T} pixels
 - $a_{zoom-in}$, $a_{zoom-out}$ **zooming** transformations of the camera by $\delta_z \%$
 - a_{null} no transformation
- we aim to learn $h_{w_h}(x_i^{(t)}, q_i) \in \mathcal{A}$

$$\boldsymbol{W}_{h} = argmin_{\boldsymbol{W}'h} \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(f_{\boldsymbol{W}}(\boldsymbol{x}_{i}^{(N_{T})}, \boldsymbol{q}_{i}), \boldsymbol{t}_{i})$$

The setup (2/3)

- Directly learning \boldsymbol{W}_h intractable
- Using Reinforcement learning to maximize the reward of an agent that controls the virtual camera
- Optimization objective: increase the probability of the VQA model answering correctly

$$r_t = \left[f_{\boldsymbol{W}} \left(\boldsymbol{x}_i^{(t)}, \boldsymbol{q}_i \right) \right]_c - \left[f_{\boldsymbol{W}} \left(\boldsymbol{x}_i^{(t-1)}, \boldsymbol{q}_i \right) \right]_c$$

 $[f_W(\cdot)]_c$ confidence of the correct answer

- Agent network: image analysis, question encoding, similarity
- Similarity from image and question extract an attention-like feature map

 $\boldsymbol{u} \in \mathbb{R}^{H_a \times W_a}$ $[\boldsymbol{u}]_{i,j} = [\widetilde{\boldsymbol{x}}]_{i,j}^T (\boldsymbol{W}_T \boldsymbol{q}) \in \mathbb{R}'$

The setup (3/3)



Reward Calculation

The frame selection problem

- Agent rewarded for every confidence increasing action to the correct answer
- Side effect: sub-optimal final action
- To overcome this: select the answer with highest average confidence
- For N_T control steps

$$\arg \max_{j} \left(\left[\frac{1}{N_{T}} \sum_{i=1}^{N_{T}} f_{W} \left(\boldsymbol{x}_{i}^{(t)}, \boldsymbol{q}_{i} \right) \right]_{j} \right)$$

Experimental setup

- VQA Model: MUTAN
- Agent:
 - Image analysis: ResNet-50 (pre-trained)
 - Question representation: GRU-based encoder (pre-trained)
 - Q-Learning training:
 - Rainbow method
 - 300,000 steps
 - replay memory of 100,000 steps
 - $\gamma = 0.99$
 - Adam optimizer, learning rate: 0.5×10^{-4}
- 5,000 episodes test from the VQA 2.0 validation set

Experimental results & Conclusions

- Active perception approach to VQA
- Trained using reinforcement learning
- transformations on the input images
- Method capable of increasing the accuracy of VQA
- High potential: frame selection can lead to significant further accuracy improvements



Question: Are there leaves on the trees? - Answer: No

Method	Accuracy	Acc. Gain
Baseline	60.36	-
Proposed (Confident Frame)	59.81	-0.55
Proposed	60.86	0.5
Proposed (Best Frame)	66.68	6.32

Acknowledgments

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 871449 (OpenDR). This publications reflects the author's views only. The European Commissions is not responsible for any use that may be made of the information it contains.



www.opendr.eu

Thank you!