

Differential Attention for Visual Question Answering

Paper: <https://arxiv.org/abs/1804.00298>

Code: https://github.com/chirag26495/DAN_VQA

Team - 13:

- Chirag Parikh - 2022900005
- Neeraj Veerla - 2021121008
- Adhiraj Deshmukh - 2021121012
- Shreya Patil - 2021121009

Overview

Visual Question Answering (The Problem)

Visual Question Answering (VQA) is a task in Machine Learning that aims to develop models capable of answering questions based on visual information.

The goal is to train a model capable of extracting features from both the image and the question, then combine them to generate an answer.

VQA is a challenging task as it requires a combination of computer vision and natural language processing techniques, since the model needs to be able to:

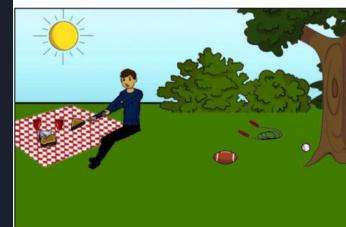
- Understand the content of the image
- Understand the meaning of the question
- Combine the visual and textual information
- Output an appropriate answer



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?



VQA-v2 Dataset

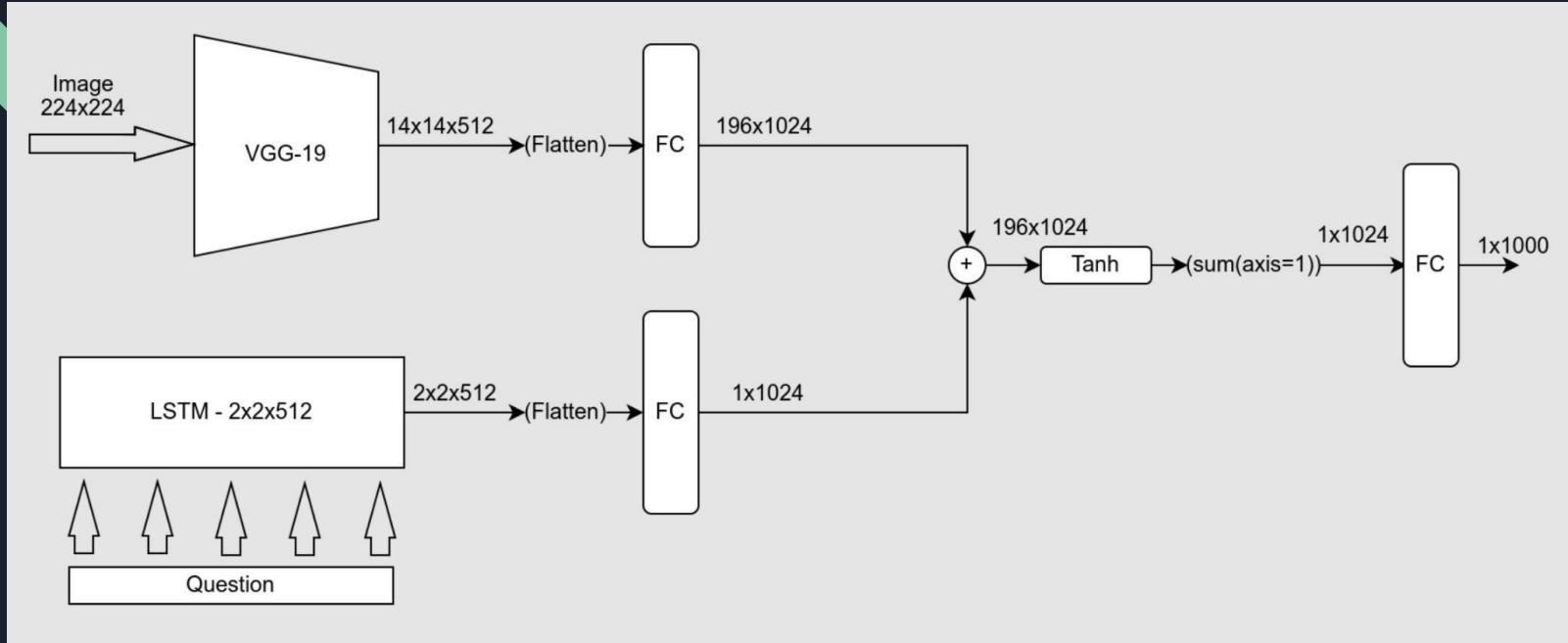
- Contains open-ended questions about images which require an understanding of vision, language and common sense knowledge to answer.
- 265,016 images (COCO and abstract scenes)
- At least 3 questions (5.4 questions on average) per image
- 10 ground truth answers per question
- **Total dataset size: 14GB**
- **Training Split: 4,43,757 (I+Q pairs)**
- **Validation Split: 2,14,354 (I+Q pairs)**



Key Components

- Basic VQA model [LSTM Q + CNN I]
- Finding Exemplars [KD-Tree and k-means clustering]
- Baseline: Simple Attention model [LSTM Q + CNN I + Attention]
- **Solution (Proposed)**: Differential Attention model [DAN + LSTM Q + CNN I + Attention]
- Other: Stacked Attention Network [SAN-2]

Basic VQA model (LSTM Q + CNN I)



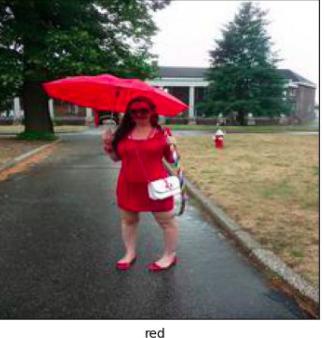
- **Joint image-question level embedding** from the above model is used to relate meaningful exemplars.
- Output: Train features of $(4,43,757 \times 1024)$ and Val features of $(2,14,354 \times 1024)$

Finding Supporting and Opposing Exemplars

Proposed method in Paper:

- For each data point in the embedding space, find the closest neighbours using KD-tree.
- Once the neighbours are found, perform clustering of the neighbours.
- Choose the nearest clusters as supporting and farthest clusters as opposing.
- The main idea was not choosing opposing examples too far from the data point.

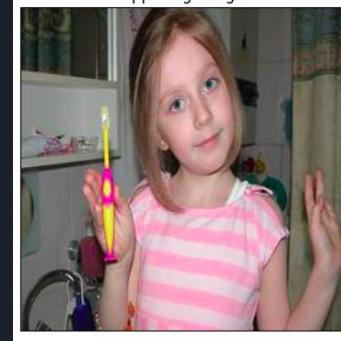
What color is her umbrella?



Supporting image



Opposing image



Finding Supporting and Opposing Exemplars (cont.)

- The method proposed by the author is computationally expensive.
Estimated runtime of 59 days on ADA, with time complexity of
 $O(D \cdot [N \cdot \log|D| + N \cdot k \cdot t])$
- Alternative method:
 - We choose the nearest 50 data-points as supporting exemplars,
 - 1200-1500 nearest neighbours for opposing exemplars.
 - Reduces time complexity to $O(N \cdot D \cdot \log|D|)$
 - Reduces the time of computation by a factor of $\approx 10^2$
 - With multiprocessing, reduced the time to ~ 12 hrs

Time for computation per instance:

Method	Query-time	k-means	compute-distance	Total
KNN + K-Means	0.7231 s	11.8403 s	0.0002 s	12.75 s
KNN	0.583 s	-	0.0002 s	0.591 s

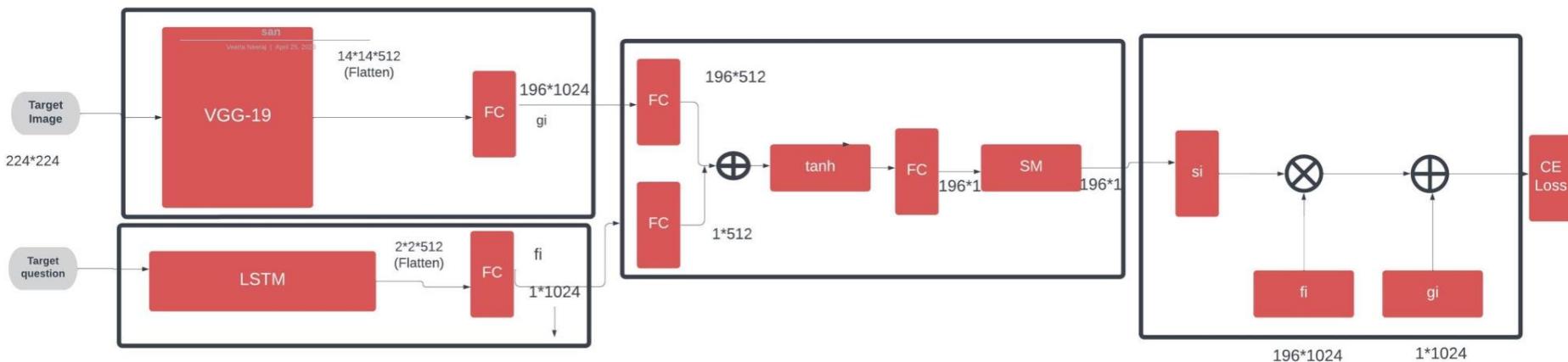
Finding Supporting and Opposing Exemplars (cont.)

Our Alternative method:

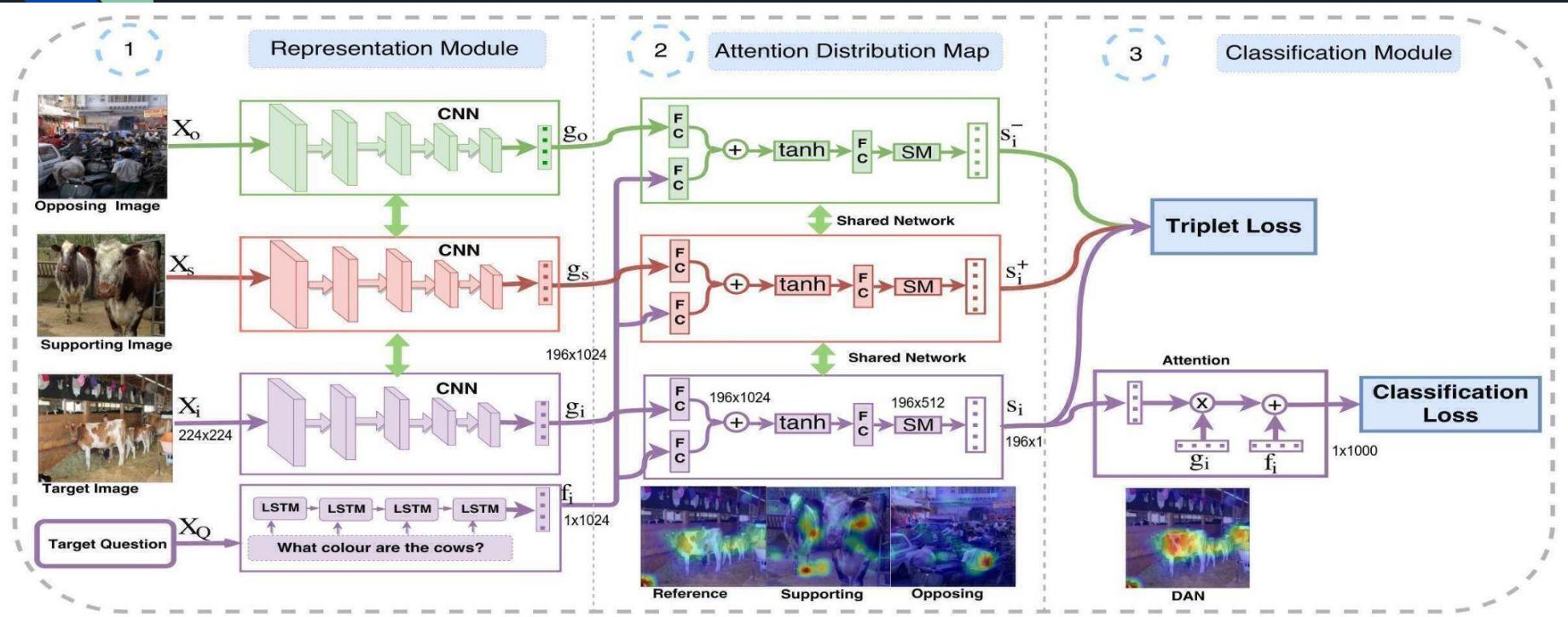
- We first clustered the whole data using Mini Batch K-means using multiprocessing.
- Performed KD-Tree on the cluster centres of the data-points.
- For all the points in a particular cluster, we choose the nearest clusters to be the supporting and the clusters in range 20 to 25 to be opposing clusters.
- This method was tried out to reduce the time complexity of the find examples in just 10 mins.



Baseline: (LSTM Q + CNN I + Attention)



Proposed: Differential Attention Network (DAN) + Baseline





Triplet Loss

- Alpha margin (=0.2) was empirically chosen, on the basis of loss on validation dataset.
- Si, Si+, Si- are the attention weighted regions of target image, supporting exemplar and opposing exemplar.

$$\begin{aligned} T(s_i, s_i^+, s_i^-) \\ = \max(0, \|t(s_i) - t(s_i^+)\|_2^2 + \alpha - \|t(s_i) - t(s_i^-)\|_2^2) \end{aligned}$$

Cross-Entropy Loss

- C (=1000) is no. of answer categories.
- v (=1) is the weightage of Triplet Loss in the Total loss.

$$\begin{aligned} L(\mathbf{s}, \mathbf{y}, \theta) &= \frac{1}{N} \sum_{i=1}^N (L_{cross}(\mathbf{s}, \mathbf{y}) + \nu T(s_i, s_i^+, s_i^-)) \\ L_{cross}(\mathbf{s}, \mathbf{y}) &= -\frac{1}{C} \sum_{j=1}^C y_j \log p(c_j | \mathbf{s}) \end{aligned}$$



Experiments

Different approaches for Finding Exemplars: (Author's and our Alternative)

- Different sets of exemplars from two different methods. (1-50/75/100, 1000-1200/1500)

Variants of DAN architecture:

- Different outputs from vgg-net. (7x7, 14x14, freeze/unfreeze)
- Different channel sizes for FC layers. (1024, 512)

We tried different input feature types for the triplet loss:

- Attention features (ha)
- Weighted attention scores (pi)
- Attended image features (vi_attended)
- Updated question features (u)



Triplet Loss hyperparameters tuning

Margin:

- The margin helps in defining a decision boundary for the embeddings in the feature space.
- The model learns to keep the positive examples within a certain distance from the anchor example and the negative examples beyond a certain distance from the anchor example.
- This ensures that the embeddings of examples of the same class are closer together and those of different classes are farther apart in the feature space.

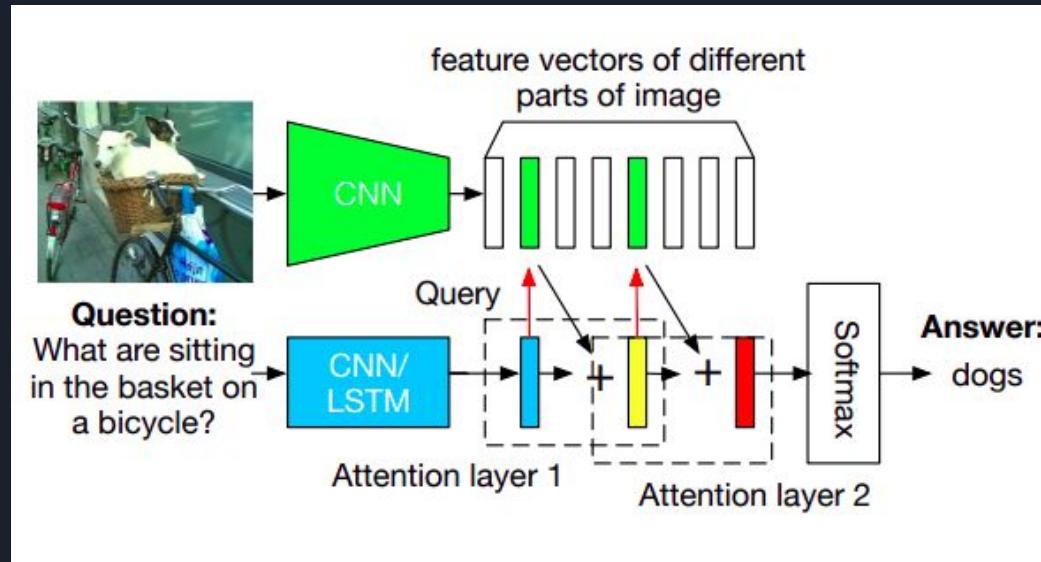
We tried to increase the margin from 0.2 to 0.8, **0.2** works the best, this is also dependent on how we find the exemplars and the parameters used in the exemplar finding.

V Weightage:

- Weighs the importance of triplet loss as compared to the CE loss, in the final loss of the model.

From 1 to 10, **10** works fine, as triplet loss decrease rapidly and the value is small.

Other model: Stacked Attention Networks (SAN-2)



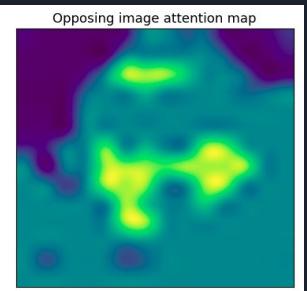
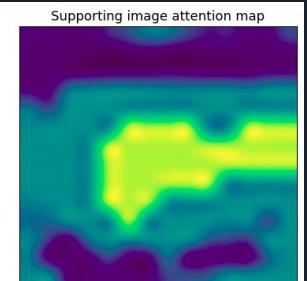
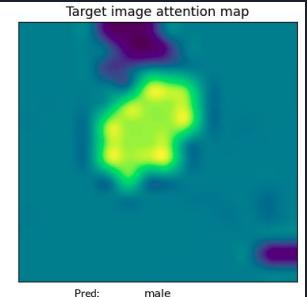
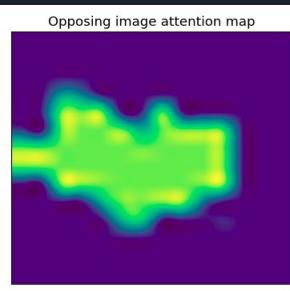
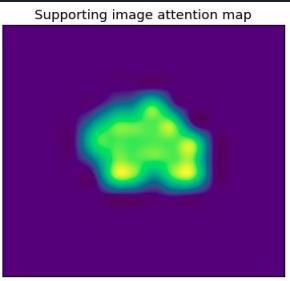
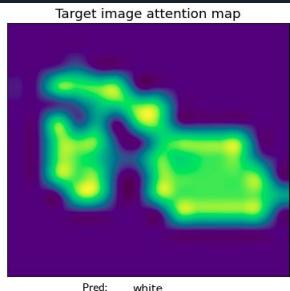
- Stacking an additional Attention layer over our Baseline model (using updated question features) for SAN-2.

Quantitative Results

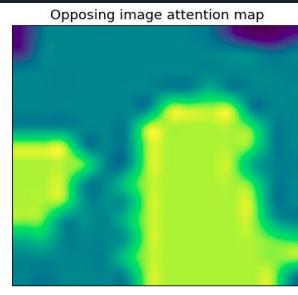
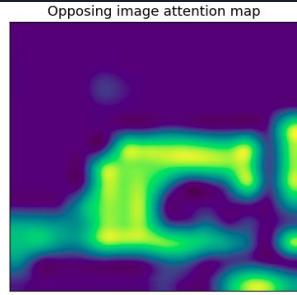
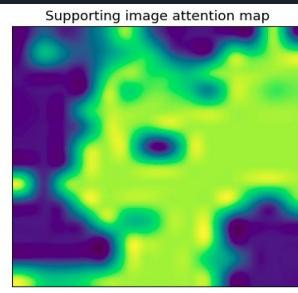
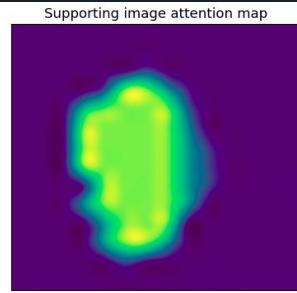
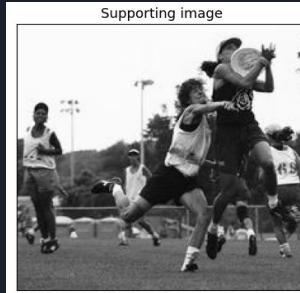
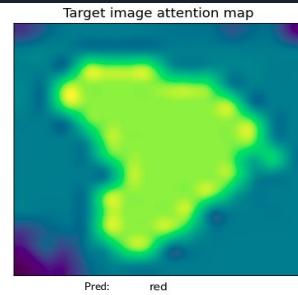
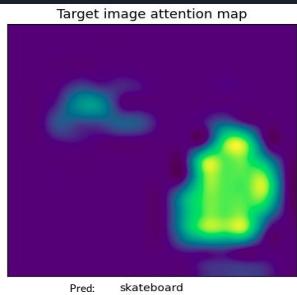
VQA2.0 accuracy on Validation set

<u>Models</u>	<u>All</u>	<u>Yes/No</u>	<u>Number</u>	<u>Other</u>
Basic (LQI)	47.61	74.86	37.21	29.61
Baseline (LQIA)	53.23	76.39	38.53	39.49
SAN-2	55.28	77.15	39.69	42.76
DAN + LQIA	55.49	77.23	39.59	42.55
DAN-alt. + LQIA	54.16	77.13	38.75	40.77

Qualitative Results of DAN



Qualitative Results of DAN (contd.)



Thank you!

Prof - Avinash Sharma

TA's -

- Amogh
- Pratyush

