# Deep Learning for Vision & Language
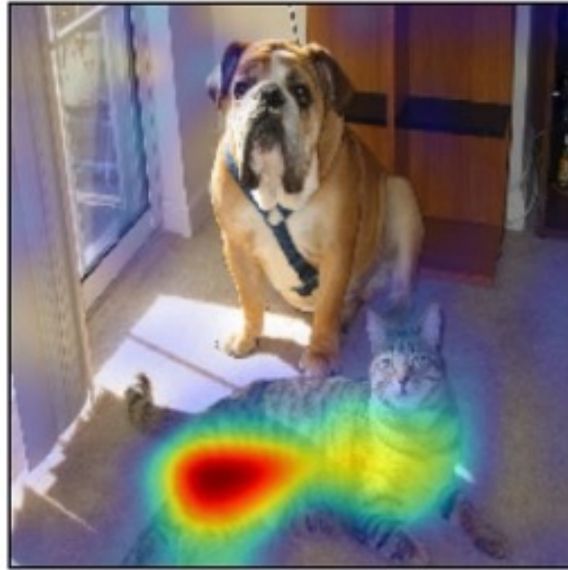
Other Topics

RICE UNIVERSITY

# Today

- GradCAM: Explaining Predictions
- ALBEF
- ALBEF + AMC
- Two-step Optimization for Object Detection
- Visual Question Answering: General Framework
- Mobile UI Navigation: Spotlight

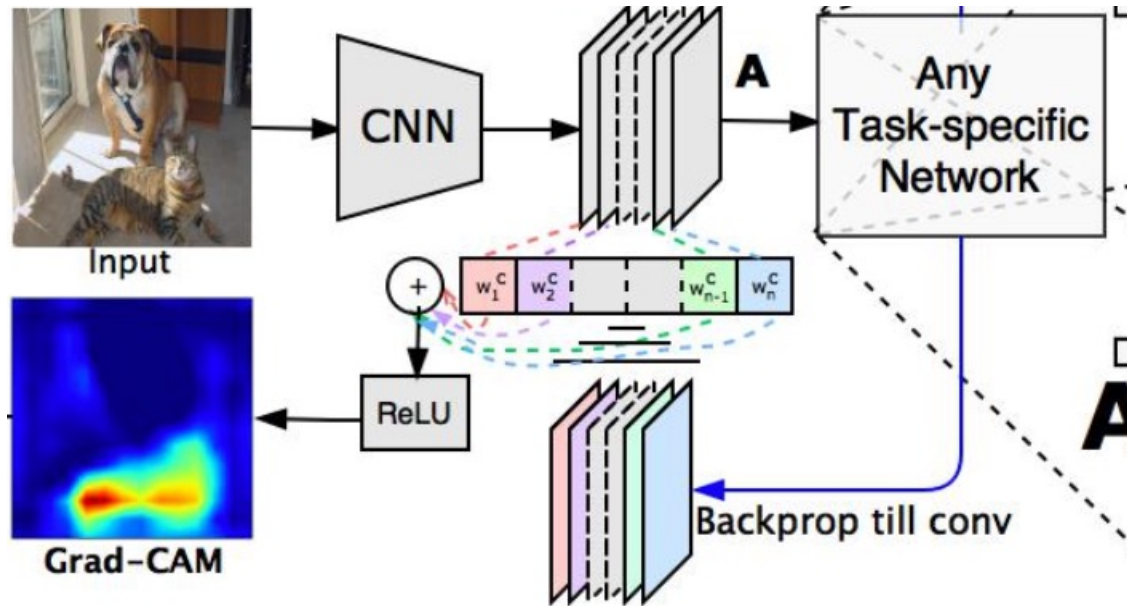# Explainability: GradCAM



(a) Original Image  (c) Grad-CAM 'Cat'  (i) Grad-CAM 'Dog'

https://arxiv.org/abs/1610.02391

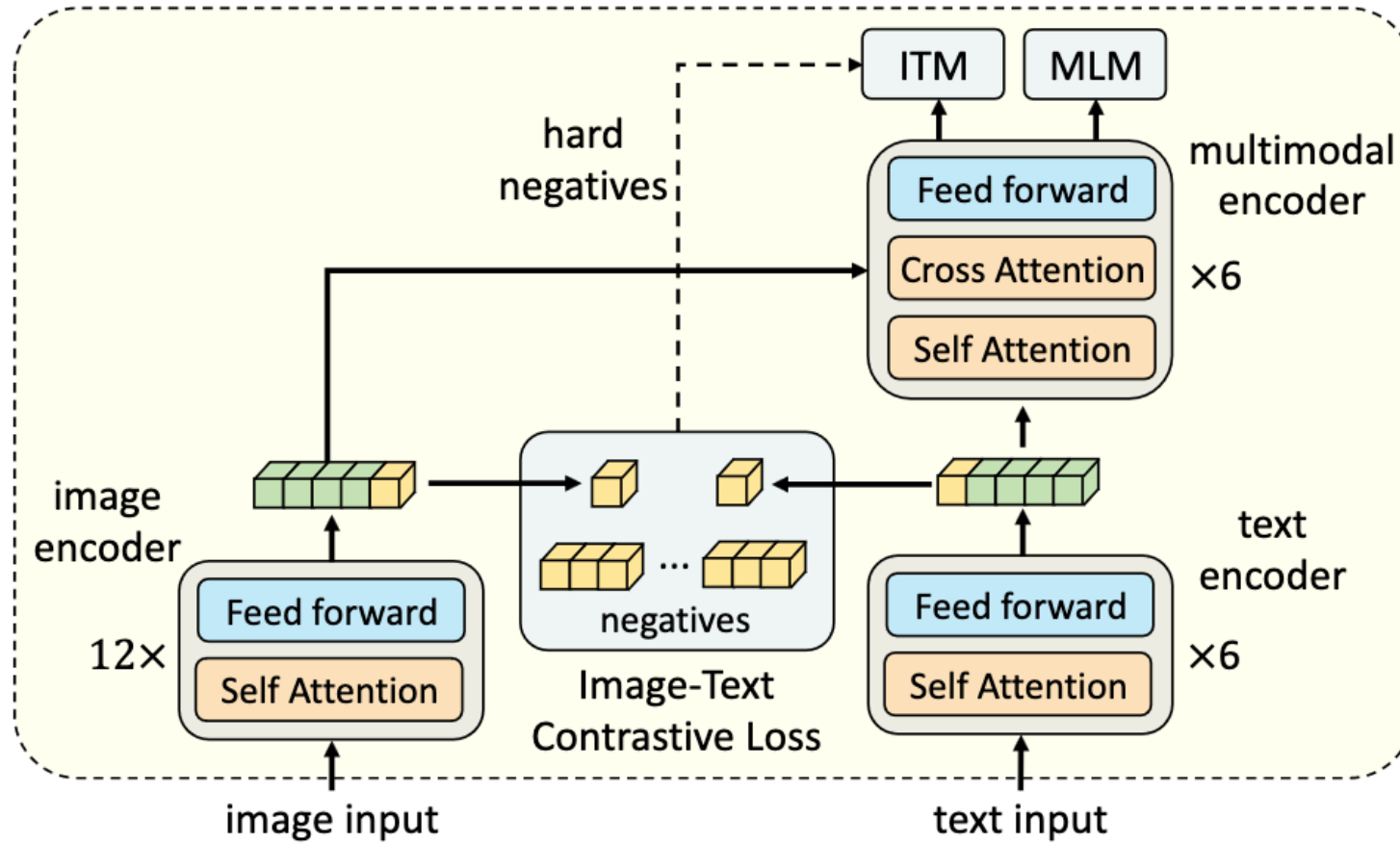# Explainability: GradCAM (2017)



$$\alpha_k^c = \overbrace{\frac{1}{Z} \sum_i \sum_j}^{\text{global average pooling}} \underbrace{\frac{\partial y^c}{\partial A_{ij}^k}}_{\text{gradients via backprop}}$$
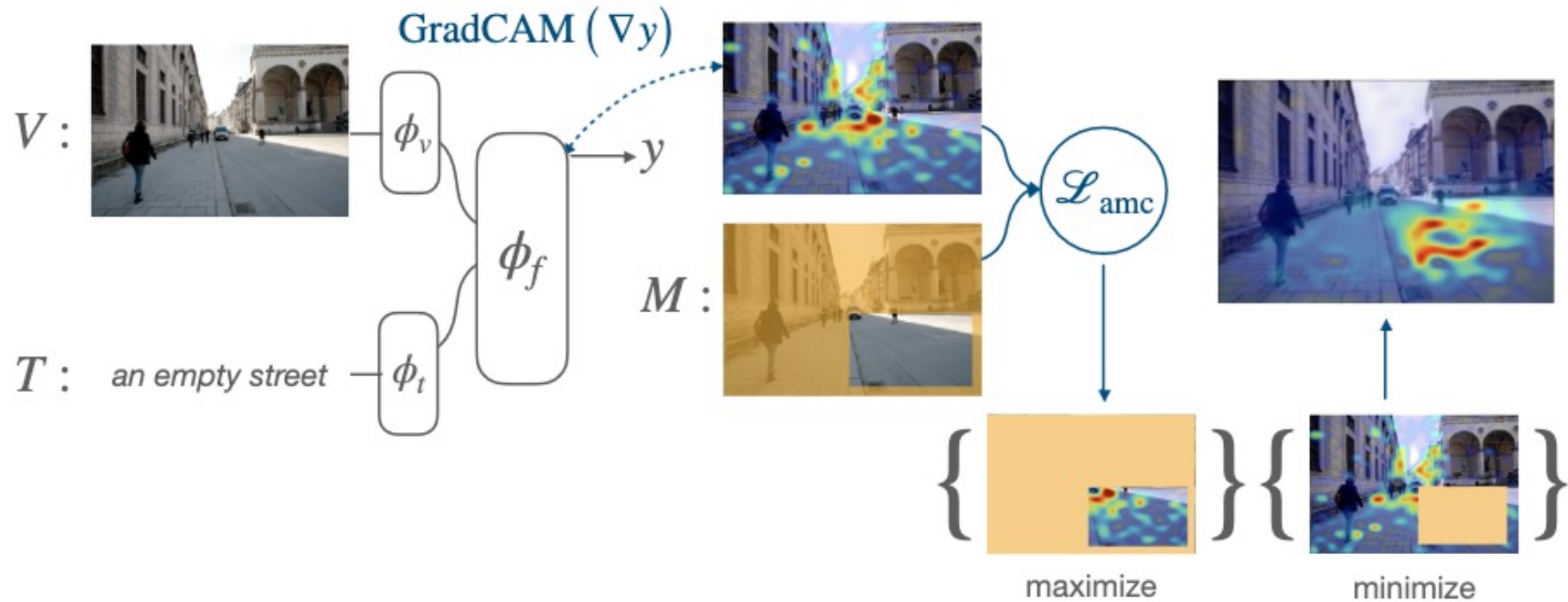
$$L_{\text{Grad-CAM}}^c = ReLU \underbrace{\left( \sum_k \alpha_k^c A^k \right)}_{\text{linear combination}}$$

# Explainability with Vision-Language Models
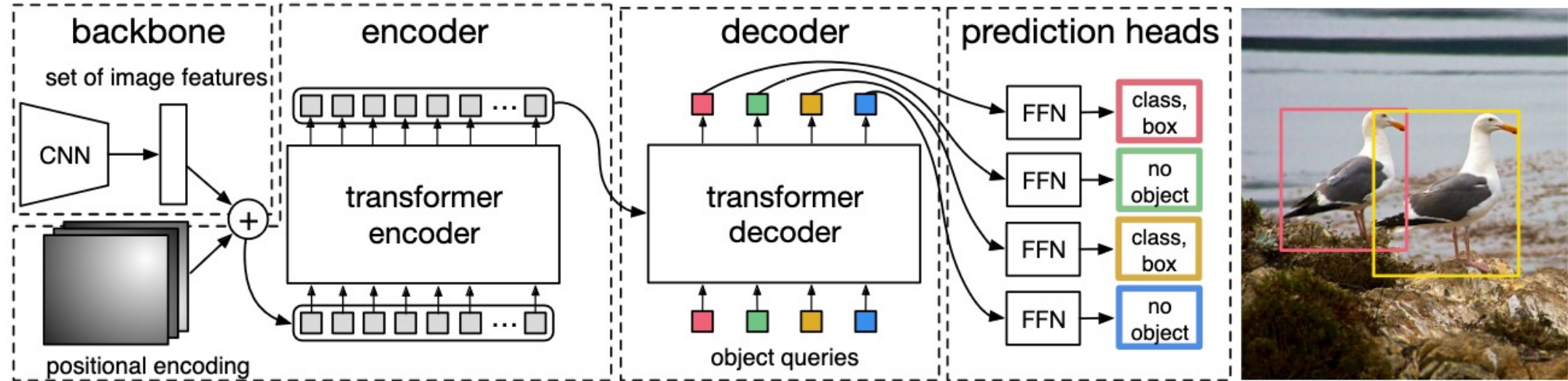# Case Study: The ALBEF model

# Attention Mask Consistency (AMC)



$$\mathcal{L}_{\max} = \mathbb{E}_{(V,T,M)\sim D} \max\left(0, \max_{i,j}\left((1 - M_{i,j})\,A_{i,j}\right) - \max_{i,j}\left(M_{i,j}A_{i,j}\right) + \Delta_2\right)$$

https://arxiv.org/abs/2206.15462

# Object Detection with Transfromers (DETR) (2020)



$$\mathcal{L}_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^{N} \left[ -\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{\{c_i \neq \varnothing\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}}(i)) \right]$$

where $\hat{\sigma} = \arg\min_{\sigma \in \mathfrak{S}_N} \sum_{i}^{N} \mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)})$

# VQA: Visual Question Answering

www.visualqa.org

Aishwarya Agrawal*, Jiasen Lu*, Stanislaw Antol*,
Margaret Mitchell, C. Lawrence Zitnick, Dhruv Batra, Devi Parikh

| Is this person trying to hit a ball? | yes<br>yes<br>yes | yes<br>yes<br>yes |
| --- | --- | --- |
| What is the person hitting the ball with? | frisbie<br>racket<br>round paddle | bat<br>bat<br>racket |

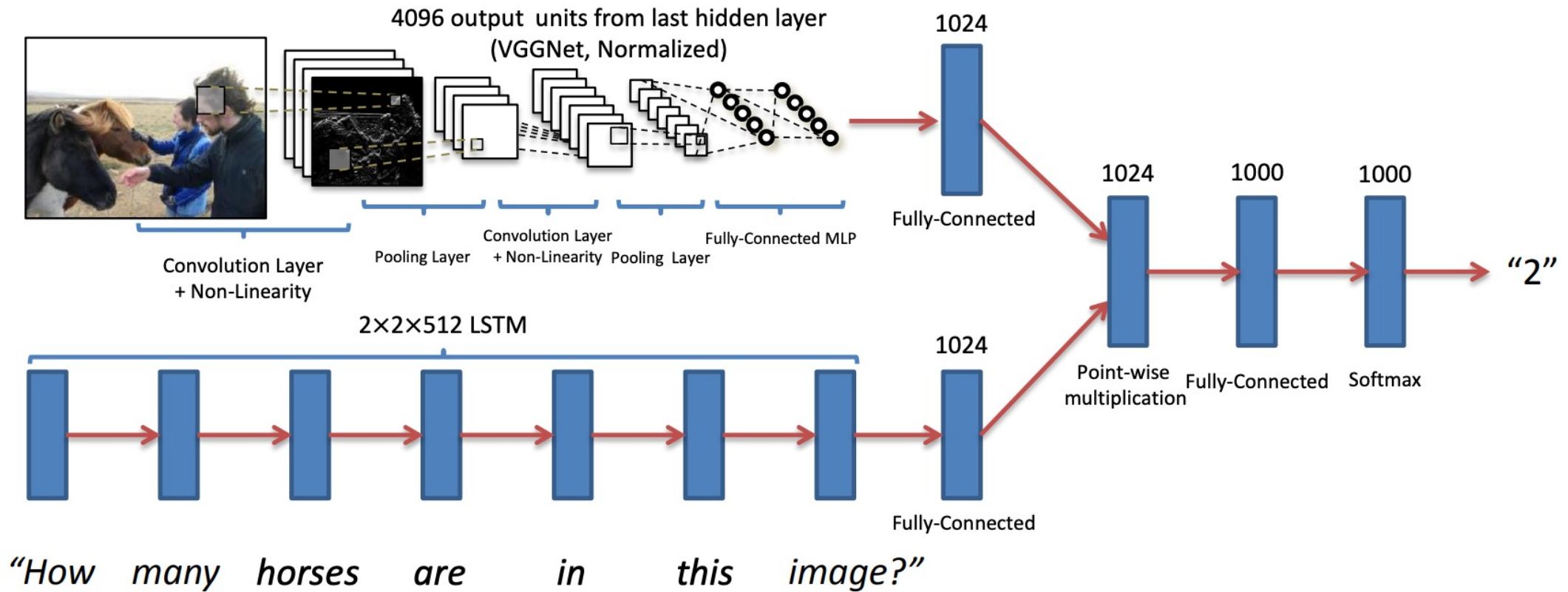| What is the guy doing as he sits on the bench? | phone<br>taking picture<br>taking picture with phone | reading<br>reading<br>smokes |
| --- | --- | --- |
| What color are his shoes? | blue<br>blue<br>blue | black<br>black<br>brown |

# Visually Grounded Question Answering

https://arxiv.org/pdf/1505.00468.pdf

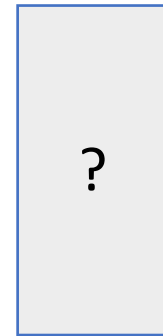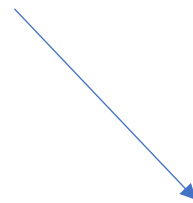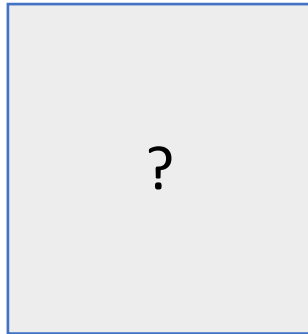# What Features to use as input visual features?

CVPR 2017

# In Defense of Grid Features for Visual Question Answering

Huaizu Jiang[1,2]*, Ishan Misra[2], Marcus Rohrbach[2], Erik Learned-Miller[1], and Xinlei Chen[2]
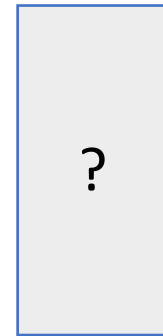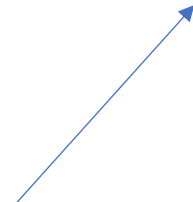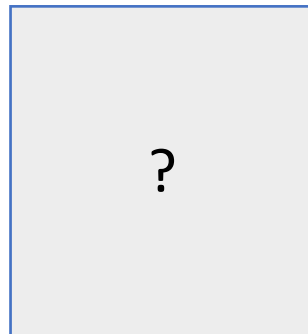
[1]UMass Amherst, [2]Facebook AI Research (FAIR)

{hzjiang,elm}@cs.umass.edu, {imisra,mrf,xinleic}@fb.com

# VQA Solution today?



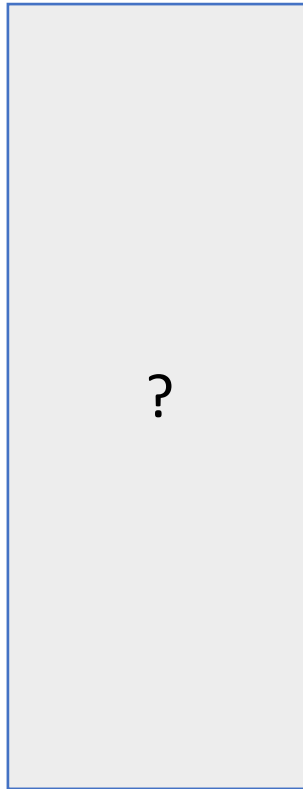What is the color of the jacket of the man on this picture?

?

?

?

Cross Entroy Loss Across 5000 possible answers

# VQA Solution today?



What is the color of the jacket of the man on this picture?
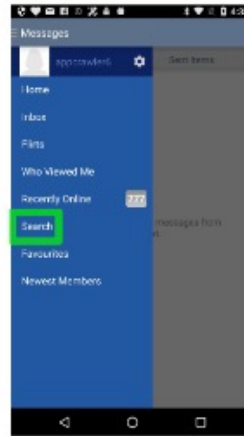
?

Cross Entroy Loss Across 5000 possible answers

# Spotlight: Visual Interface Navigation



**Command Grounding**
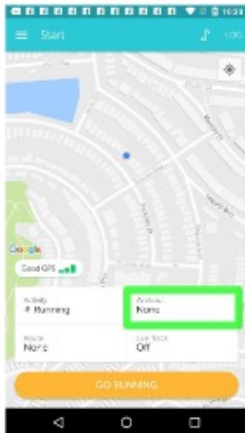
Select zoom in button | Yes
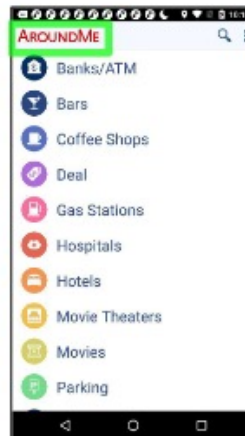
Select icon above home | No

click on the second image from the bottom second row | Yes
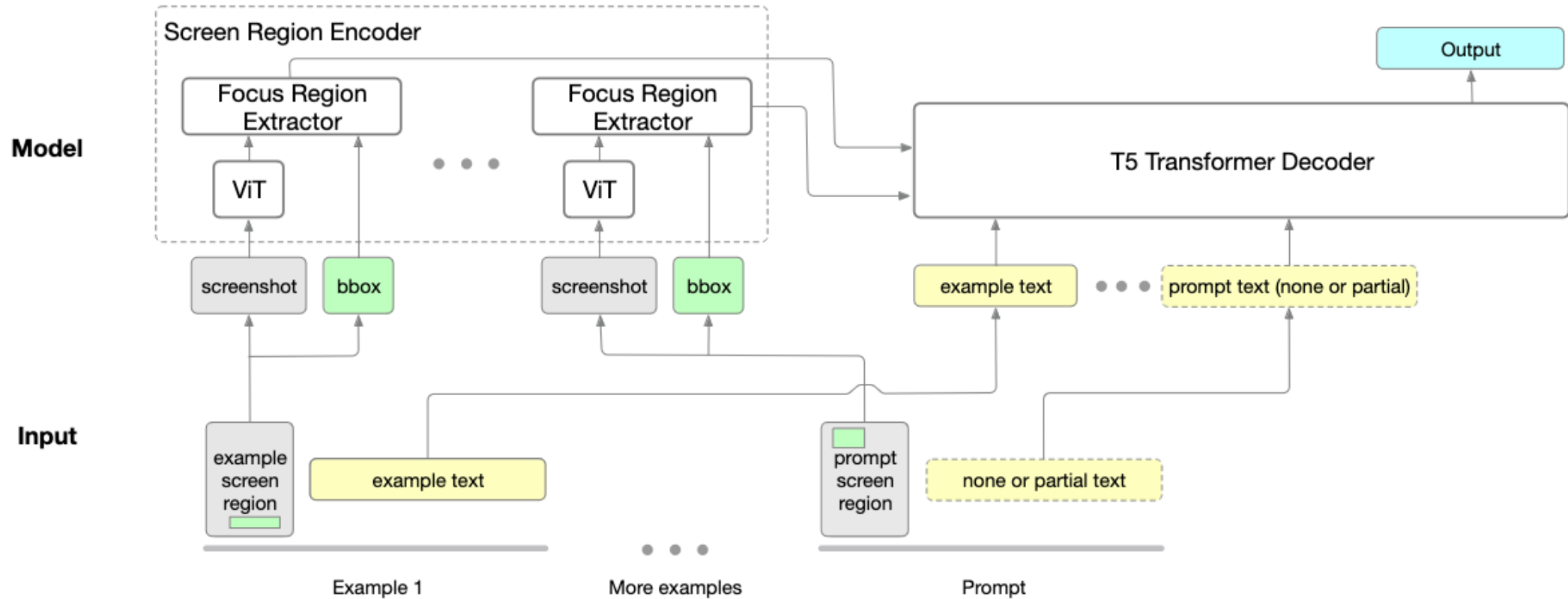
**Tappability Prediction**

Tappable | Yes

Tappable | No

Tappable | Yes

# Spotlight: Visual Interface Navigation

# Questions