# NEURO-SYMBOLIC VISUAL REASONING:
## DISENTANGLING "VISUAL" FROM "REASONING"

**SAEED AMIZADEH**
SAAMIZAD@MICROSOFT.COM

**ALEX POLOZOV**
POLOZOV@MICROSOFT.COM

**HAMID PALANGI**
HPALANGI@MICROSOFT.COM

**YICHEN HUANG**
YICHUANG@MIT.EDU

**KAZUHITO KOISHIDA**
KAZUKOI@MICROSOFT.COM

# VISUAL QUESTION ANSWERING

[**GQA:** Hudson & Manning, 2019]

**Q**: *"What color is the food on the red object left of the small girl that is holding a hamburger?"*

**VQA Model**

**A**: *"Yellow."*

**Language Signal**

**Reasoning**

**Visual Perception**

**Answer**

**Visual Signal**

# REASONING ≙ LOGICAL REASONING + EXTRA CAPABILITIES

Pure logical reasoning does not often suffice for visual reasoning because visual perception is noisy and uncertain.



$x$

**Example:** imperfect visual perception classifies $\Pr(\text{Husky} \mid x) \approx \Pr(\text{Wolf} \mid x)$.

Then,
$\mathbf{Pr}$("*Is there a **husky** in the living room?*") ≈
$\mathbf{Pr}$("*Is there a **wolf** in the living room?*")

Yet "*in the living room*" or the visual context should resolve the ambiguity.

# RESEARCH QUESTIONS

1. Given a visual featurization $\mathcal{V}$ of a visual scene, how informative $\mathcal{V}$ is on its own to answer a question about the scene without learned reasoning?

2. How solvable is VQA/GQA given perfect vision?

3. For an arbitrary VQA model $\mathcal{M}$, how much its reasoning abilities can compensate for the imperfections in perception to solve the task?

# OUR CONTRIBUTIONS



**(I) Differentiable First-Order Logic ($\nabla$-FOL) for Visual Description & Reasoning**



**(II) Evaluation of Reasoning vs. Perception for VQA models using $\nabla$-FOL**

# FIRST ORDER LOGIC FOR SCENE DESCRIPTION

**Scene Graph Representation**



**FOL Representation**

"<u>There is</u> a <u>cat</u> to the <u>left</u> of <u>all</u> objects."

$$F_Q : F(X, Y) = \exists X \forall Y : Cat(X) \wedge Left(X, Y)$$

- **Variables** enumerates over detected objects.

- **Atomic Predicates** represent object names, attributes and binary relations.

- **Formulas** represent a statement or a question about the scene.

# FOL FOR POSING A HYPOTHETICAL QUESTION

**Scene Graph Representation**



**FOL Representation**

"There is a cat to the left of all objects."

$$F_Q: F(X,Y) = \exists X \forall Y: Cat(X) \wedge Left(X,Y)$$

"Is there a cat to the left of all objects?"

This question can be answered **probabilistically** by evaluating the likelihood:

$$\alpha(F_Q) \triangleq \Pr(Answer = "Yes"|I) = \Pr(F_Q \Leftrightarrow True|I)$$

**exponentially hard to calculate directly** ☹

# ∇-FOL: INFERENCE IN POLYNOMIAL TIME

In order to do inference in polynomial time, we introduce the intermediate notion of attention on the object $x_i$ w.r.t. formula $F$:

$$\alpha(F|x_i) \triangleq \mathbf{Pr}\big(F_{X=x_i} \Leftrightarrow True\big), \quad \textbf{Where} \quad F_{X=x_i} \triangleq F(x_i, Y, \dots, Z)$$

Then the answer likelihood can be reduced to computing attention via aggregation operators $A_\forall$ and $A_\exists$:

If $X$ is universally quantified ($\forall$):

$$\alpha(F) = \prod_{i=1}^{N} \alpha(F|x_i) \triangleq A_\forall(\alpha(F|X))$$

If $X$ is existentially quantified ($\exists$):

$$\alpha(F) = 1 - \prod_{i=1}^{N} (1 - \alpha(F|x_i)) \triangleq A_\exists(\alpha(F|X))$$

# ∇-FOL: RECURSIVE CALCULATION OF ATTENTION



**Negation Operator**
$$\alpha(F|x_i) = 1 - \alpha(G|x_i) \triangleq \mathbf{Neg}[\alpha(G|x_i)]$$

**Filter Operator**
$$\alpha(F|x_i) = \alpha(\pi|x_i) \cdot \alpha(G|x_i) \triangleq \mathbf{Filter}_\pi[\alpha(G|x_i)]$$

**Relate Operator**
$$\alpha(F|x_i) = A_q\left(\left[\bigodot_{\pi\in\Pi_{XY}} \alpha(\pi|x_i, Y)\right] \odot \alpha(G|Y)\right)$$
$$\triangleq \mathbf{Relate}_{q,\Pi_{XY}}[\alpha(G|Y)],$$
$$\forall i \in [1..N], q = Quantifier(Y) \in \{\exists, \forall\}$$

**Every FOL formula**

**NOT** **Smaller FOL formula**

**Unary Predicate** **AND** **Smaller FOL formula**

**Binary Predicate** **AND** **Smaller FOL formula**

8/14/2020

# THE LANGUAGE SYSTEM: FROM NATURAL LANGUAGE TO FOL FORMULA

**Natural Language**

"Is there a **ball** **on** the **table**?"

*Semantic parsing*

**Task-dependent DSL**

Select (**Table**) → Relate(**on**, **Ball**) → Exists(?)

*Compilation*

**Task-independent ∇-FOL**

$A_\exists(\text{Filter}_{\textbf{Ball}}[\text{Relate}_{\textbf{on},\exists}[\text{Filter}_{\textbf{Table}}[1]]])$

*Equivalence*

**First-order Logic**

$\exists X, \exists Y: \textbf{Ball}(X) \wedge \textbf{Table}(Y) \wedge \textbf{On}(X,Y)$

NEURO-SYMBOLIC VISUAL REASONING

# GQA DOMAIN SPECIFIC LANGUAGE

| GQA OP | T | Equivalent FOL Description | Equivalent DFOL Program |
|---|---|---|---|
| $\mathbf{GSelect}(name)[]$ | N | $name(X)$ | $\mathbf{Filter}_{name}[1]$ |
| $\mathbf{GFilter}(attr)[\alpha_X]$ | N | $attr(X)$ | $\mathbf{Filter}_{attr}[\alpha_X]$ |
| $\mathbf{GRelate}(name, rel)[\alpha_X]$ | N | $name(Y) \wedge rel(X, Y)$ | $\mathbf{Relate}_{rel,\exists}\big[\mathbf{Filter}_{name}[\alpha_X]\big]$ |
| $\mathbf{GVerifyAttr}(attr)[\alpha_X]$ | Y | $\exists X : attr(X)$ | $\mathcal{A}_\exists(\mathbf{Filter}_{attr}[\alpha_X])$ |
| $\mathbf{GVerifyRel}(name, rel)[\alpha_X]$ | Y | $\exists Y \exists X : name(Y) \wedge rel(X, Y)$ | $\mathcal{A}_\exists(\mathbf{Relate}_{rel,\exists}\big[\mathbf{Filter}_{name}[\alpha_X]\big])$ |
| $\mathbf{GQuery}(category)[\alpha_X]$ | Y | $[\exists X : c(X) \text{ for } c \text{ in } category]$ | $[\mathcal{A}_\exists(\mathbf{Filter}_c[\alpha_X]) \text{ for } c \text{ in } category]$ |
| $\mathbf{GChooseAttr}(a_1, a_2)[\alpha_X]$ | Y | $[\exists X : a(X) \text{ for } a \text{ in } [a_1, a_2]]$ | $[\mathcal{A}_\exists(\mathbf{Filter}_a[\alpha_X]) \text{ for } a \text{ in } [a_1, a_2]]$ |
| $\mathbf{GChooseRel}(n, r_1, r_2)[\alpha_X]$ | Y | $[\exists Y \exists X : n(Y) \wedge r(X, Y) \text{ for } r \text{ in } [r_1, r_2]]$ | $[\mathcal{A}_\exists(\mathbf{Relate}_{r,\exists}\big[\mathbf{Filter}_n[\alpha_X]\big]) \text{ for } r \text{ in } [r_1, r_2]]$ |
| $\mathbf{GExists}()[\alpha_X]$ | Y | $\exists X...$ | $\mathcal{A}_\exists(\alpha_X)$ |
| $\mathbf{GAnd}()[\alpha_X, \alpha_Y]$ | Y | $\exists X... \wedge \exists Y...$ | $\mathcal{A}_\exists(\alpha_X) \cdot \mathcal{A}_\exists(\alpha_Y)$ |
| $\mathbf{GOr}()[\alpha_X, \alpha_Y]$ | Y | $\exists X... \vee \exists Y...$ | $1 - (1 - \mathcal{A}_\exists(\alpha_X)) \cdot (1 - \mathcal{A}_\exists(\alpha_Y))$ |
| $\mathbf{GTwoSame}(category)[\alpha_X, \alpha_Y]$ | Y | $\exists X \exists Y \bigvee_{c \in category} (c(X) \wedge c(Y))$ | $\mathcal{A}_\exists([\mathcal{A}_\exists(\mathbf{Filter}_c[\alpha_X]) \cdot \mathcal{A}_\exists(\mathbf{Filter}_c[\alpha_Y]) \text{ for } c \text{ in } category])$ |
| $\mathbf{GTwoDifferent}(category)[\alpha_X, \alpha_Y]$ | Y | $\exists X \exists Y \bigwedge_{c \in category} (\neg c(X) \vee \neg c(Y))$ | $1 - \mathbf{GTwoSame}(category)[\alpha_X, \alpha_Y]$ |
| $\mathbf{GAllSame}(category)[\alpha_X]$ | Y | $\bigvee_{c \in category} \forall X : c(X)$ | $1 - \prod_{c \in category} (1 - \mathcal{A}_\forall(\mathbf{Filter}_c[\alpha_X]))$ |

# VISUAL SYSTEM: FROM IMAGE TO PREDICATES



$$\alpha("cat"|x_i)$$
$$\alpha("Dog"|x_i)$$
$$\alpha("Man"|x_i)$$
.
.
.

**Off-the-shelf Object Detection (e.g. Faster-RCNN, Ren et al. 2015)**

$x_i$

Object Detection

Object Featurization

Neural Visual Oracle

**Neural Visual Oracle**

Cat    Dog    Man    • • •

**Queried Predicates**

# THE WHOLE SYSTEM



Q: "Is there a man on the left of all objects in the scene?"

A: "No"

Semantic Parser

$$F(X,Y) = \exists X: Man(X) \wedge \forall Y: Left(X,Y)$$

Differentiable FOL Reasoning

Visual Featurization

$u_1$
$u_2$
$u_3$
$u_4$

$X$ {
1   $\alpha(man|x_1)$
1   $\alpha(man|x_2)$
1   $\alpha(man|x_3)$
1   $\alpha(man|x_4)$

Filter

Relate

1 1 1 ... 1

$\alpha(F|x_1)$
$\alpha(F|x_2)$
$\alpha(F|x_3)$
$\alpha(F|x_4)$

Aggregation

$\alpha(F)$

Loss Function

$\alpha(man|x_i)$    $\alpha(Left|x_i,y_j)$   $\alpha(F|y_1) \dots \alpha(F|y_4)$

Visual Oracle

# USING $\nabla$-FOL TO EVALUATE PERCEPTION

**Q1:** Given a visual featurization $\mathcal{V}$ for a certain VR task, how informative $\mathcal{V}$ is on its own to solve the task using mere FOL for reasoning?

**For GQA:** The visual featurization $\mathcal{V}$ is the Faster-RCNN featurization [Ren et. al, 2015].

# BUILDING THE BASE MODEL

| The Base Model |
|---|

1) Put $\nabla$-FOL on the top of a neural Visual Oracle $\mathcal{O}$.

2) Train the resulted architecture using the Faster-RCNN featurization, the golden programs and golden answers in GQA via indirect supervision from the answer.

3) Denote the result as the *Base Model $\mathcal{M}_\phi$*.

# USING $\nabla$-FOL TO EVALUATE PERCEPTION

**Q1:** Given a <span style="color:red">visual featurization</span> $\mathcal{V}$ for a certain VR task, <span style="color:red">how informative</span> $\mathcal{V}$ is on its own to solve the task using mere FOL for reasoning?

| Split | Accuracy | Consistency |
|-------|----------|-------------|
| Open | 42.73 % | 88.74 % |
| Binary | 65.08 % | 86.65 % |
| All | 51.86 % | 88.35 % |

Table 1: The accuracy and consistency on Test-Dev for the Base model using the Faster-RCNN features.

$\nabla$-FOL has <span style="color:red">no trainable parameters</span>, so the **accuracy** of $\mathcal{M}_\phi$ on test data indirectly captures the amount of information in $\mathcal{V}$.

# USING ∇-FOL TO MEASURE THE IMPORTANCE OF PERCEPTION

**Q2**: how well a VR task can be achieved given <span style="color:red">perfect vision</span>?

**For GQA:** What happens if we replace the visual system by the <span style="color:red">Golden Scene Graphs</span>?

# BUILDING THE PERFECT MODEL

**The Perfect Model**

1) Replace the trained $\mathcal{O}$ in $\mathcal{M}_\phi$ with the golden GQA scene graphs, denoted as $\mathcal{O}^*$.
2) Denote the result as the *Perfect Model $\mathcal{M}^*$*.

# USING $\nabla$-FOL TO MEASURE THE IMPORTANCE OF PERCEPTION

**Q2**: how well a VR task can be achieved given perfect vision?

The **accuracy** of $\mathcal{M}^*$ on the GQA validation set is $\approx$ **96%**.

Achieving such high upper-bound shows that:

➤The $\nabla$-FOL is sound.

➤The GQA task is heavily vision-dependent.

# USING ∇-FOL TO EVALUATE REASONING

**Q3:** How much the <span style="color:red">reasoning abilities</span> of a candidate model $\mathcal{M}$ can <span style="color:red">compensate</span> for the <span style="color:red">imperfections in perception</span> to solve the task?

**Important:** $\mathcal{M}$ is arbitrary! Need not be DFOL-based.

**For GQA:** we compare MAC Network [Hudson & Manning, 2018] vs LXMERT [Tan & Bansal, 2019].

# HARD SET VS EASY SET



Test-Dev

Base Model $\mathcal{M}_\phi$

Easy Set   Hard Set

The **accuracy** of $\mathcal{M}$ on the hard set $(\mathbf{Acc}_h)$ captures the amount the reasoning process of $\mathcal{M}$ compensates for its imperfect perception.

The **error** of $\mathcal{M}$ on the easy set $(\mathbf{Err}_e)$ captures the degree to which the reasoning process of $\mathcal{M}$ distorts the informative visual signals.

# USING ∇-FOL TO EVALUATE REASONING

**Q3:** How much the reasoning abilities of a candidate model $\mathcal{M}$ can compensate for the imperfections in perception to solve the task?

| | Split | Test-Dev | | Hard Test-Dev | | Easy Test-Dev | |
|---|---|---|---|---|---|---|---|
| | | Accuracy | Consistency | $Acc_h$ | Consistency | $Err_e$ | Consistency |
| MAC | Open | 41.66 % | 82.28 % | 18.12 % | 74.87 % | 26.70 % | 84.54 % |
| | Binary | 71.70 % | 70.69 % | 58.77 % | 66.51 % | 21.36 % | 75.37 % |
| | All | 55.37 % | 79.13 % | 30.54 % | 71.04 % | 23.70 % | 82.83 % |
| LXMERT | Open | **47.02 %** | **86.93 %** | **25.27 %** | **85.21 %** | **22.92 %** | **87.75 %** |
| | Binary | **77.63 %** | **77.48 %** | **63.02 %** | **73.58 %** | **13.93 %** | **81.63 %** |
| | All | **61.07 %** | **84.48 %** | **38.43 %** | **81.05 %** | **17.87 %** | **86.52 %** |

Table 2: The $\mathbf{Acc}_h$, $\mathbf{Err}_e$ and consistency for MAC and LXMERT over Test-Dev and its hard and easy subsets according to the Base model.

# CONCLUSION REMARKS

In this work, we

1. Proposed a differentiable visual description and reasoning formalism directly derived from first order logic.

2. Proposed coherent methodology for separately evaluating perception and reasoning using our differentiable first order logic formalism.

3. Incorporated our framework for the **GQA** task and two of its famous models and arrived at insightful observations.

## Thank you ☺

# SUPPLEMENTAL MATERIALS

# MODELING OPEN QUESTIONS USING FOL

For open questions, we generate all potential options for the answer, treat each option as a binary question and choose the one with highest likelihood.

For example: "*What is the color of the ball on the left of all objects*?" can be answered by answering a set of binary questions:

"Is the ball on the left of all objects blue?" → $\mathbf{Pr}(F_{Q_1} \Leftrightarrow True|I)$

"Is the ball on the left of all objects red?" → $\mathbf{Pr}(F_{Q_2} \Leftrightarrow True|I)$

"Is the ball on the left of all objects green?" → $\mathbf{Pr}(F_{Q_3} \Leftrightarrow True|I)$

# BEYOND PURE LOGICAL REASONING:
# TOP-DOWN CONTEXTUAL CALIBRATION



**Example of a reasoning technique beyond pure DFOL:**

Reminder: suppose $\alpha("Husky"|x) \approx \alpha("Wolf"|x)$.

Then, $\mathbf{Pr}("Is\ there\ a\ husky\ in\ the\ living\ room?\ ") \approx \mathbf{Pr}("Is\ there\ a\ wolf\ in\ the\ living\ room?\ ")$

However, the context "*in the living room*" should help resolve the ambiguity.

In other words, the context can be used to **calibrate** the attentions values in the top-down manner.

# BEYOND PURE LOGICAL REASONING: TOP-DOWN CONTEXTUAL CALIBRATION

Instead of uniform, assume the attention values $\alpha(F|x)$ are Beta distributed, then the posterior is:

$$\mathbf{Pr}(F \Leftrightarrow \top \,|\, \alpha) = \frac{c\alpha^w}{c\alpha^w + d(1-c)(1-\alpha)^v}$$

Where $c, d, w, v$ are derived from the beta distribution likelihood + the prior and are estimated from the question context using a bi-LSTM.

# EFFECT OF TOP-DOWN CONTEXTUAL CALIBRATION

| | | Test-Dev | | Hard Test-Dev | | Easy Test-Dev | |
|---|---|---|---|---|---|---|---|
| | Split | Accuracy | Consistency | $\text{Acc}_h$ | Consistency | $\text{Err}_e$ | Consistency |
| $\nabla$-FOL | Open | **41.22 %** | **87.63 %** | **0.53 %** | **11.46 %** | **2.53 %** | 90.70 % |
| | Binary | 64.65 % | **85.54 %** | 4.42 % | 61.11 % | **2.21 %** | **86.33 %** |
| | All | 51.45 % | **87.22 %** | 1.81 % | 19.44 % | **2.39 %** | **89.90 %** |
| Calibrated $\nabla$-FOL | Open | **41.22 %** | 86.37 % | **0.53 %** | **11.46 %** | **2.53 %** | 89.45 % |
| | Binary | **71.99 %** | 79.28 % | **37.82 %** | **70.90 %** | 9.20 % | 84.45 % |
| | All | **54.76 %** | 84.48 % | **12.91 %** | **57.72 %** | 6.32 % | 88.51 % |

Table 3: The $\textbf{Acc}_h$, $\textbf{Err}_e$ and consistency for $\nabla$-FOL and Calibrated $\nabla$-FOL over Test-Dev and its hard and easy subsets.