# Hard to Cheat: A Turing Test based on Answering Questions about Images

Mateusz Malinowski and Mario Fritz
{mmalinow, mfritz}@ mpi-inf.mpg.de
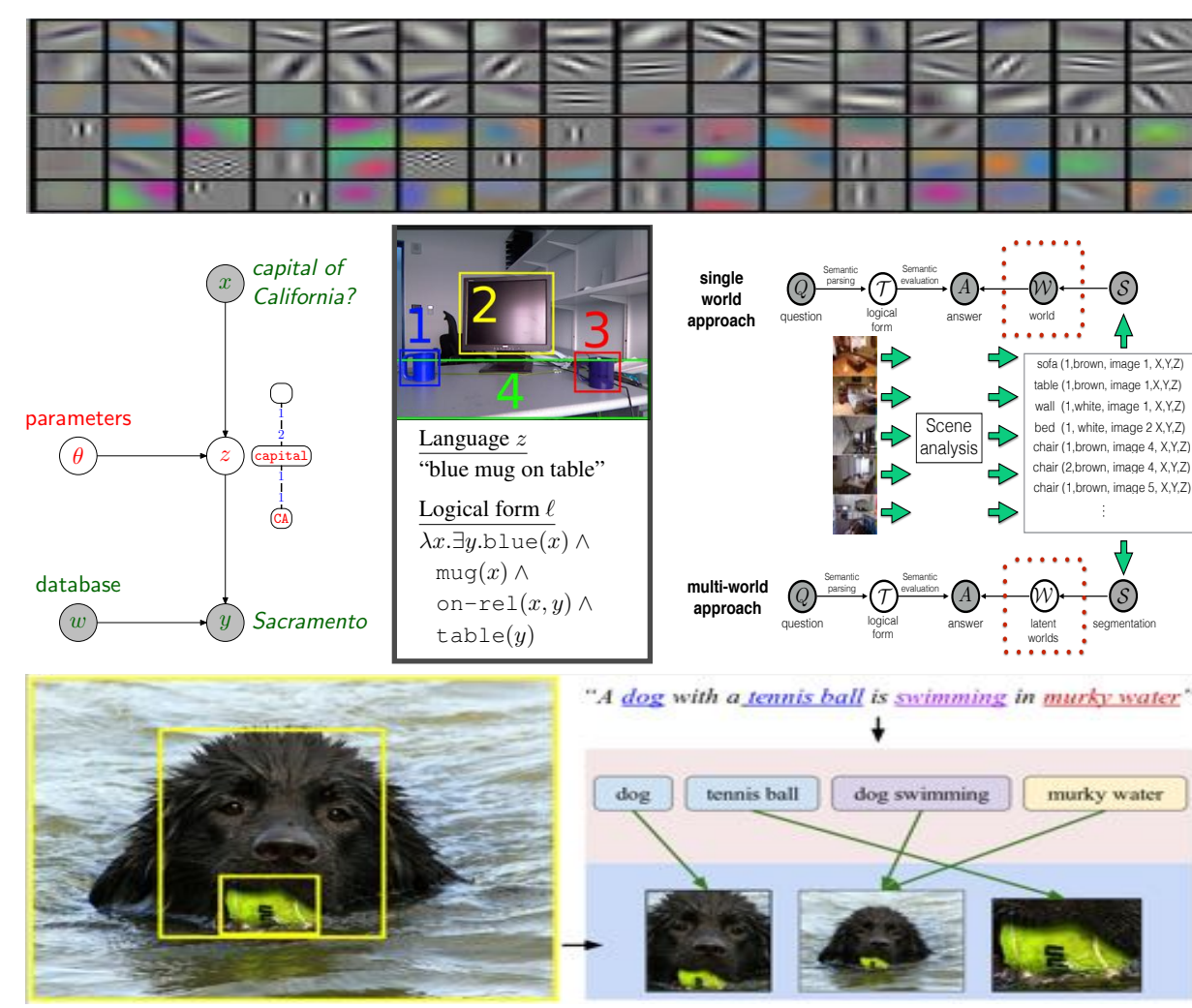www.d2.mpi-inf.mpg.de/visual-turing-challenge

## Motivation

- Stronger vision and language techniques are being developed
- Can machines answer on questions about natural images?
- A holistic, open-ended, end-to-end test that resembles the famous TT
- No internal representation is evaluated; challenge is open to diverse approaches
- Likely to be less prone to over interpretation that TT
- Scalable annotation effort
- Strategies for automatic evaluation

## Related work

- Machine perception
- Machine language understanding
- Grounding
- Image-to-sentence alignment
- Question-answering problem



## Overview

- Introduce a holistic Visual Turing Challenge
- Discuss associated challenges in Vision and NLP
- Introduce and discuss performance measures
  - Social consensus to benchmark different architectures



## Challenges

- Vision and language
  - Joint treatment of both modalities
    - 'Which hand of the teacher is on her chin?'
    - Ideally closing the loop for improved perception
  - Richness of the concepts
    - Object categories
    - Attributes (e.g. genders, colors, states)
    - Unknown human notion of spatial relations
  - Ambiguities in the reference frame
    - Object-centric
    - Observer-centric
    - World-centric
  - Contextualization of the concepts
    - White in 'white elephant' and 'white snow'
- Common sense knowledge
  - Narrows down likely options or locations
    - 'Which object on the table is used for cutting?'
    - 'What is in front of scissors?'
- Defining a benchmark
  - End-to-end system that learns from textual question-answer pairs
  - Internal representation of architectures is irrelevant
  - Easy to collect a dataset
  - Hard to define automatic performance measures

## Annotations

- Unique advantages of question answering task over other tasks in terms of acquisition and task evaluation
- Cheaper annotations as no logical forms or image annotations are required
- Methods are judged not on an internal representation but provided answers
  - The task is agnostic to internal representation of a method
- Easier to formulate evaluation due to restricted output space
  - TT and language generation tasks can be challenging to evaluate
- Harder to cheat: likely robustness to over-interpretations
  - The task requires answering to the point rather than cheating an interrogator by giving generic answers that are open to interpretations

## Metrics

- Automatic Evaluation by Design
- Ambiguity
  - Cultural bias
  - Fined grained categorization
  - Reference frame
- 'Soft' Accuracy

$$\frac{1}{N}\sum_{i=1}^{N}\min\{\prod_{a\in A^i}\max_{t\in T^i}\mu(a,t),\ \prod_{t\in T^i}\max_{a\in A^i}\mu(a,t)\}\cdot 100$$

- Coverage in the lexical databases
- Further development of the metrics
  - Consider multiple human answers
  - Interpretation metric
    - Maximal score over different human answers
  - Consensus metric
    - Average over different human answers
    - Takes an agreement between human responses into account
- Experimental scenarios
  - Controlled and open scenarios with another resources available in training

## Conclusions

- Visual Turing Test provides a rich set of challenges in Vision and NLP
- Annotation and evaluation remain tractable
- Less prone to "overinterpretation"
- Automatic benchmarking, but coverage can be an issue
- Cultural bias, changes in the reference frame, naming ambiguities, and unknown spatial relation are inherent to the challenge

## DAQUAR

- NYU-Depth V2 dataset with textual question-answer pairs
- 1449 RGBD indoor images
- About 12500 question-answer pairs
- About 9 question-answer pairs per image
- Object category occurs 4 times in training set
- Answers are: colors, numbers, objects and sets of these
- First result established in [1] with comparison to human performance
- Discussion of challenges in [2]



QA: (What is behind the table?, window)
Spatial relation like 'behind' are dependent on the reference frame. Here the annotator uses observer-centric view.

QA: (what is beneath the candle holder, decorative plate)
Some annotators use variations on spatial relations that are similar, e.g. 'beneath' is closely related to 'below'.

The annotators are using different names to call the same things. The names of the brown object near the bed include 'night stand', 'stool', and 'cabinet'.

Some objects, like the table on the left of image, are severely occluded or truncated. Yet, the annotators refer to them in the questions.

QA: (what is in front of the wall divider?, cabinet)
Annotators use additional properties to clarify object references (i.e. wall divider). Moreover, the perspective plays an important role in these spatial relations interpretations.

QA: (what is behind the table?, sofa)
Spatial relations exhibit different reference frames. Some annotators use observer-centric, others object-centric.
QA: (how many lights are on?, 6)
Moreover, some questions require detection of states 'light on or off'.

QA1:(How many doors are in the image?, 1)
QA2:(How many doors are in the image?, 5)
Different interpretation of 'door' results in different counts: 1 door at the end of the hall vs. 5 doors including lockers

QA: (How many drawers are there?, 8)
The annotators use their common-sense knowledge for amodal completion. Here the annotator infers the 8th drawer from the context

QA1: (what is in front of the curtain behind the armchair?, guitar)
QA2: (what is in front of the curtain?, guitar)
Spatial relations matter more in complex environments where reference resolution becomes more relevant. In cluttered scenes, pragmatism starts playing a more important role

Q: what is at the back side of the sofas?
Annotators use wide range spatial relations, such as 'backside' which is object-centric.

QA: (What is the object on the counter in the corner?, microwave)
References like 'corner' are difficult to resolve given current computer vision models. Yet such scene features are frequently used by humans.

QA: (How many doors are open?, 1)
Notion of states of object (like open) is not well captured by current vision techniques. Annotators use such attributes frequently for disambiguation.

[1] M. Malinowski and M. Fritz "A Multi-World Approach to Question Answering about Real-World Scenes based on Uncertain Input" NIPS 2014
[2] M. Malinowski and M. Fritz "Towards a Visual Turing Challenge" NIPS Workshop on Learning Semantics 2014