# MovieQA
## Understanding Stories in Movies through Question-Answering

Makarand
Tapaswi

Yukun
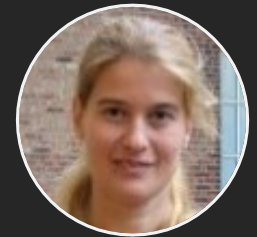Zhu

Rainer
Stiefelhagen

Antonio
Torralba

Raquel
Urtasun

Sanja
Fidler

UNIVERSITY OF TORONTO

KIT
Karlsruhe Institute of Technology

MIT
Massachusetts Institute of Technology

# MovieQA
## Understanding *stories*

| | | | |
|---|---|---|---|
| The Matrix has you | | | |

00:25:52 --> 00:25:57
Welcome, Neo. As you no doubt have guessed... I am Morpheus

00:40:42 --> 00:40:47
It exists now only as part of a neural-interactive simulation that we call the Matrix.

01:04:08 --> 01:04:09
... you know what I realize? Ignorance is bliss.

02:08:38 --> 02:08:39
Where we go from there is a choice I leave to you

**Questions:** Why does Cypher betray Morpheus? How does Trinity save Neo?

## Movie:

- 200,000 frames
- 2,000 shots
- 1,000 dialogs

- Long temporal dependencies
- Actions, interactions, emotions, intent

# MovieQA
## multiple sources of information (video and text)

**Q.** Who makes Indy return the crucifix after escaping from the grave robbers?

***A1. The local sheriff***
A2. Coronado
A3. No one, he keeps it
A4. The Boy Scout troop
A5. The grave robbers



| **PLOT** | **DVS** | **SCRIPT** | **SUBTITLE** |
|---|---|---|---|
| Indy escapes, but the local **sheriff** makes him return the **crucifix**. | Indy shows the **Cross**, more or less handing it to the **Sheriff** to make his point.<br><br>The **Sheriff** takes it casually. | SHERIFF: You still got it?<br>INDY: Well, yes sir.<br><br>Indy shows the **CROSS**, more or less handing it to the **SHERIFF** to make his point. The **Sheriff** takes it casually.<br><br>SHERIFF: I'm glad to see that | 00:10:50 --> 00:10:52<br>You still got it?<br>00:10:52 --> 00:10:53<br>Well, yes, sir.<br>00:10:55 --> 00:10:59<br>I'm glad to see that because the rightful owner of this **cross** |

# MovieQA
benchmark in numbers

- 14,944 QAs
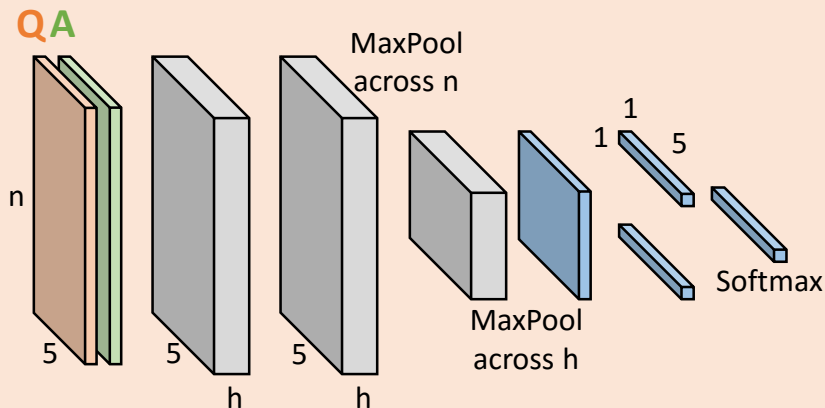- 408 movies
- 1 correct, 4 deceiving answers per Q

- 6,462 QAs with video
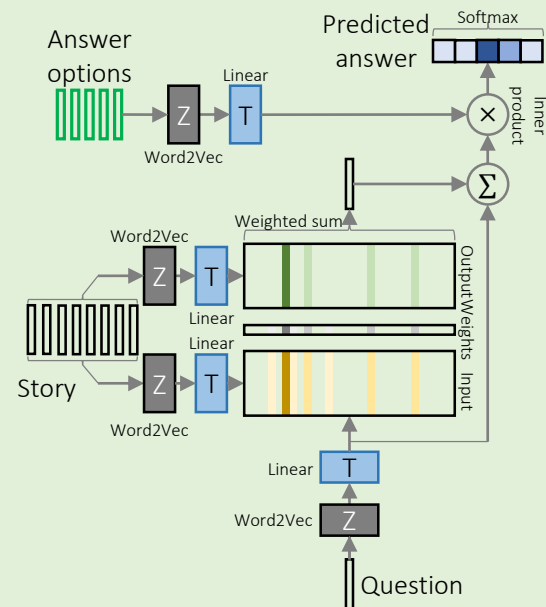- 6,771 video clips
- ~3m20s clip duration

# MovieQA
## answering methods

### General framework for multiple-choice question-answering



$$\text{correct answer} \quad \square = \arg\max_a f(\,\underset{\text{story}}{\square\square\square\square\square\square}\,,\,\underset{\text{question}}{\square}\,,\,\underset{\text{answers}}{\square\square\square\square\square}\,)$$



**Searching Student with Convolutional Brain**

QA
n
5 5 5
h h
MaxPool across n
MaxPool across h
1
1 5
Softmax



**Modified Memory Network**

Answer options
Word2Vec
Linear
Predicted answer
Softmax
Inner product
×
Σ
Weighted sum
Word2Vec
Linear
Story
Word2Vec
Linear
Output Weights
Input
Linear
Linear
Word2Vec
Question

# Video Question Answering – Video mode

- For a movie, several video "shots" are provided

- In the baseline mode, each "shot" is represented by mean-pooled representation of its frames.

# Video Question Answering – Video mode

- In this model, we propose to use attention model over shots rather than naïve mean pooling

- Feature representation based on this attention model is fed to the memory network representation, training everything together in an end-to-end fashion.