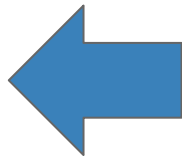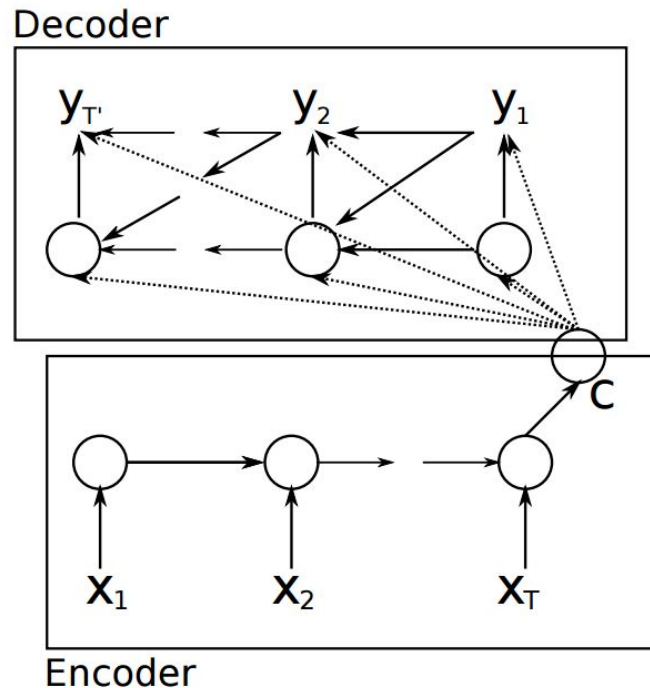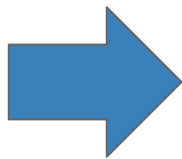# Acknowledgments



Santi Pascual
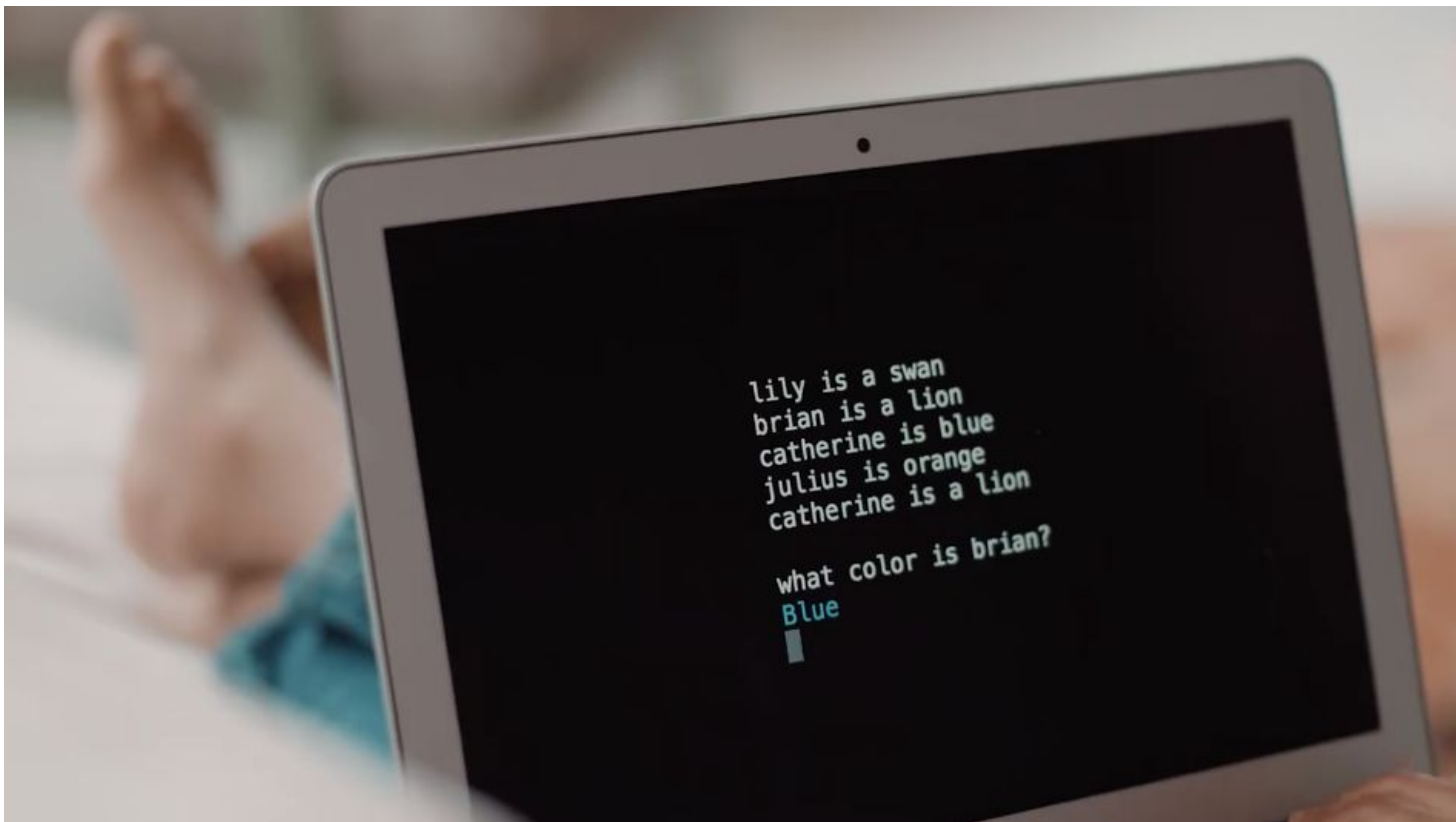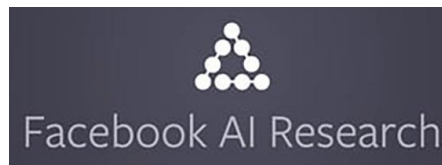
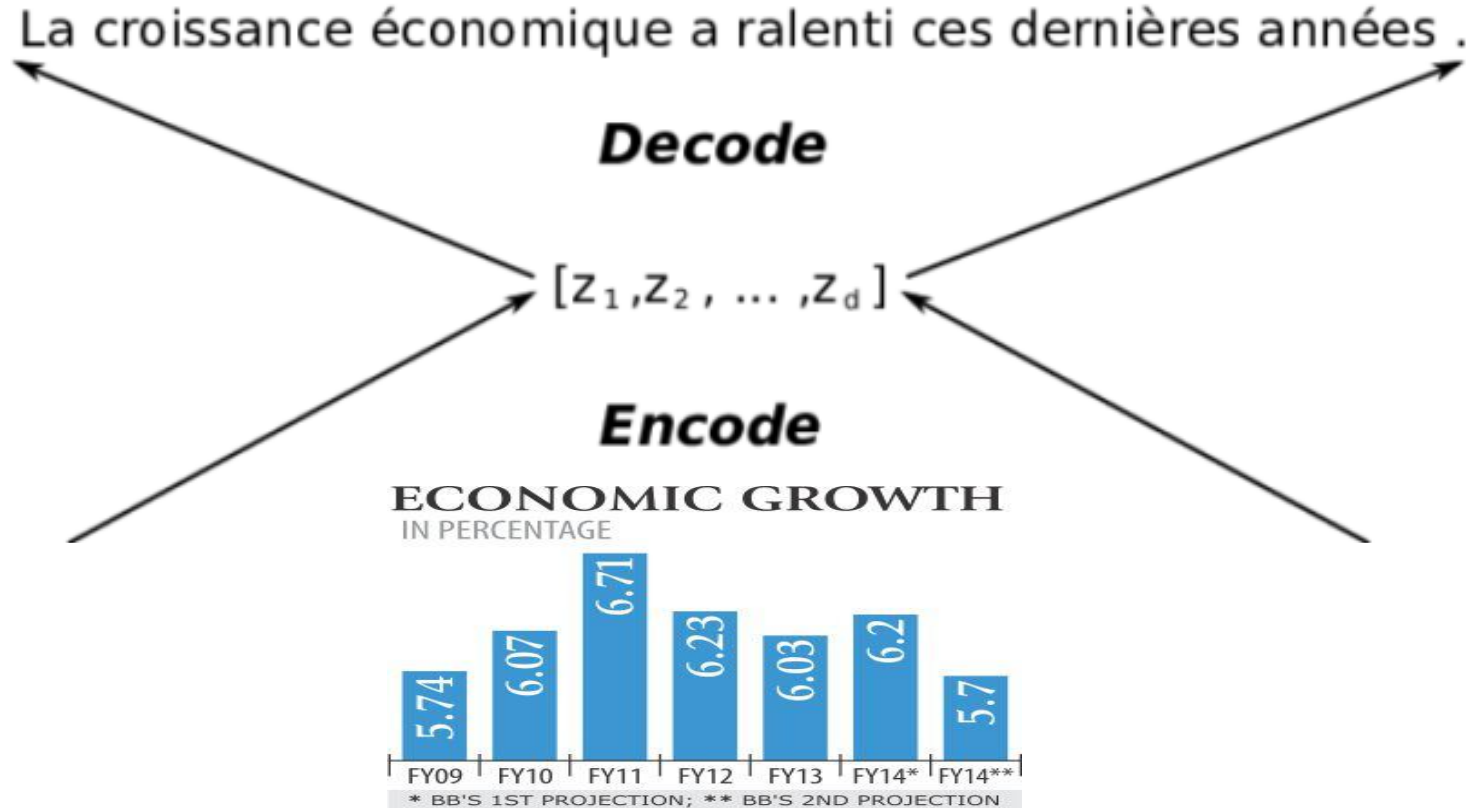# Previously in the RNN lecture...



Language OUT

Language IN

Cho, Kyunghyun, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. "Learning phrase representations using RNN encoder-decoder for statistical machine translation." arXiv preprint arXiv:1406.1078 (2014).

3

# Motivation

# Encoder-Decoder: Beyond text

La croissance économique a ralenti ces dernières années .

**Decode**

$$[z_1, z_2, \dots, z_d]$$

**Encode**

ECONOMIC GROWTH
IN PERCENTAGE

| FY09 | FY10 | FY11 | FY12 | FY13 | FY14* | FY14** |
|------|------|------|------|------|-------|--------|
| 5.74 | 6.07 | 6.71 | 6.23 | 6.03 | 6.2 | 5.7 |

* BB'S 1ST PROJECTION; ** BB'S 2ND PROJECTION

# Captioning: Show & Tell



Vinyals, Oriol, Alexander Toshev, Samy Bengio, and Dumitru Erhan. "Show and tell: A neural image caption generator." CVPR 2015.

6

# Captioning: DeepImageSent



man in black shirt is playing guitar.

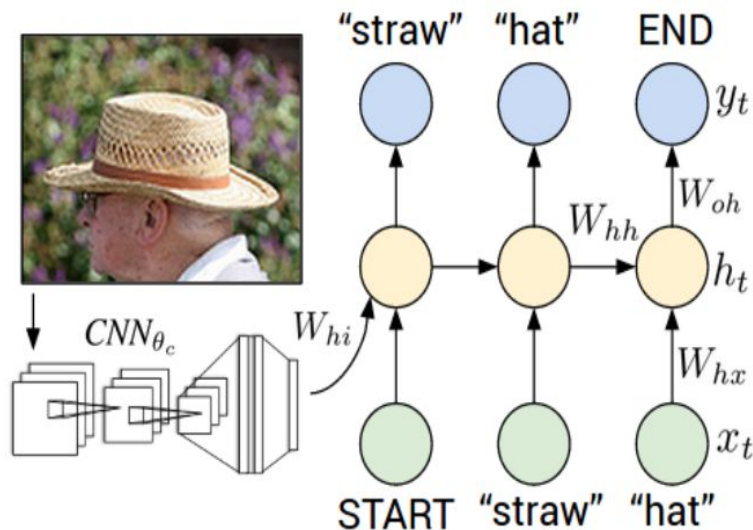construction worker in orange safety vest is working on road.

two young girls are playing with lego toy.

# Captioning: DeepImageSent

only takes into account image features in the first hidden state

$$b_v = W_{hi}[CNN_{\theta_c}(I)]$$
$$h_t = f(W_{hx}x_t + W_{hh}h_{t-1} + b_h + \boxed{\mathbb{1}(t=1) \odot b_v})$$
$$y_t = softmax(W_{oh}h_t + b_o).$$



**Multimodal Recurrent Neural Network**

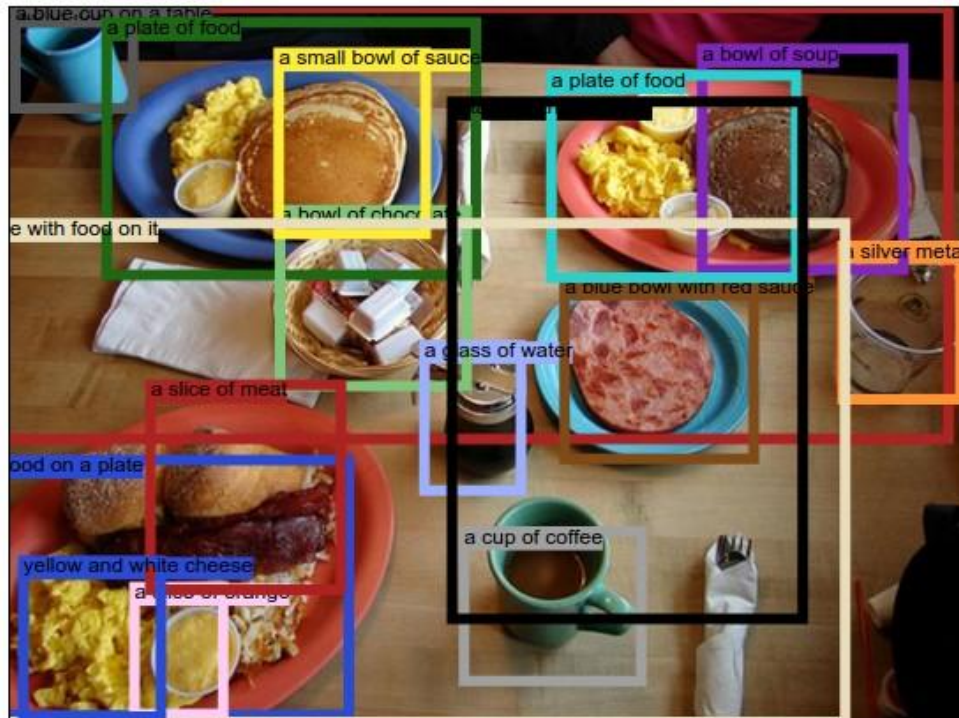# Captioning: Show & Tell



Vinyals, Oriol, Alexander Toshev, Samy Bengio, and Dumitru Erhan. "Show and tell: A neural image caption generator." CVPR 2015.

9

# Captioning (+ Detection): DenseCap



Johnson, Justin, Andrej Karpathy, and Li Fei-Fei. "Densecap: Fully convolutional localization networks for dense captioning." *CVPR 2016*

10

# Captioning (+ Detection): DenseCap



Johnson, Justin, Andrej Karpathy, and Li Fei-Fei. "Densecap: Fully convolutional localization networks for dense captioning." *CVPR 2016*

11

# Captioning (+ Detection): DenseCap



XAVI: "man has short hair", "man with short hair"

AMAIA:"a woman wearing a black shirt", "

BOTH: "two men wearing black glasses"

Johnson, Justin, Andrej Karpathy, and Li Fei-Fei. "Densecap: Fully convolutional localization networks for dense captioning." *CVPR 2016*

12

# Captioning (+ Retrieval): DenseCap



Johnson, Justin, Andrej Karpathy, and Li Fei-Fei. "Densecap: Fully convolutional localization networks for dense captioning." *CVPR 2016*
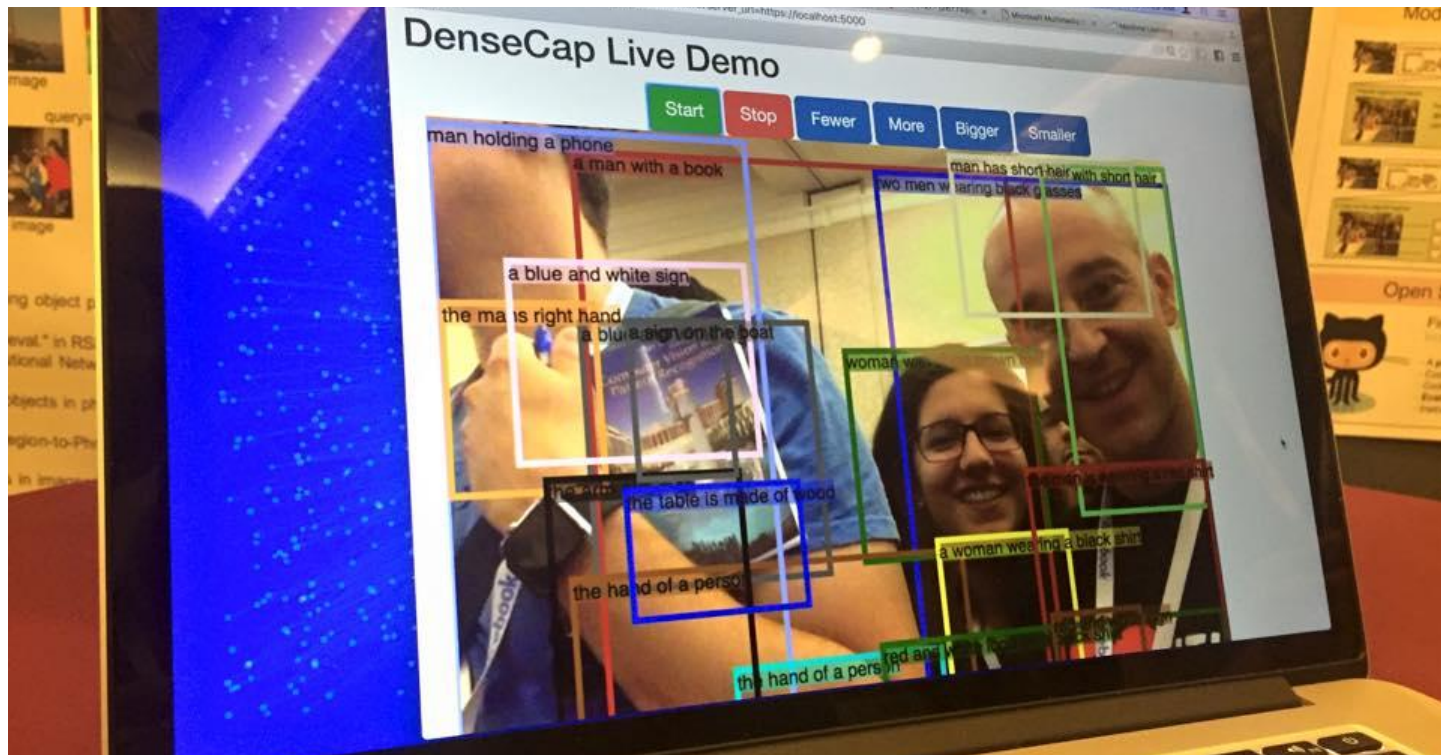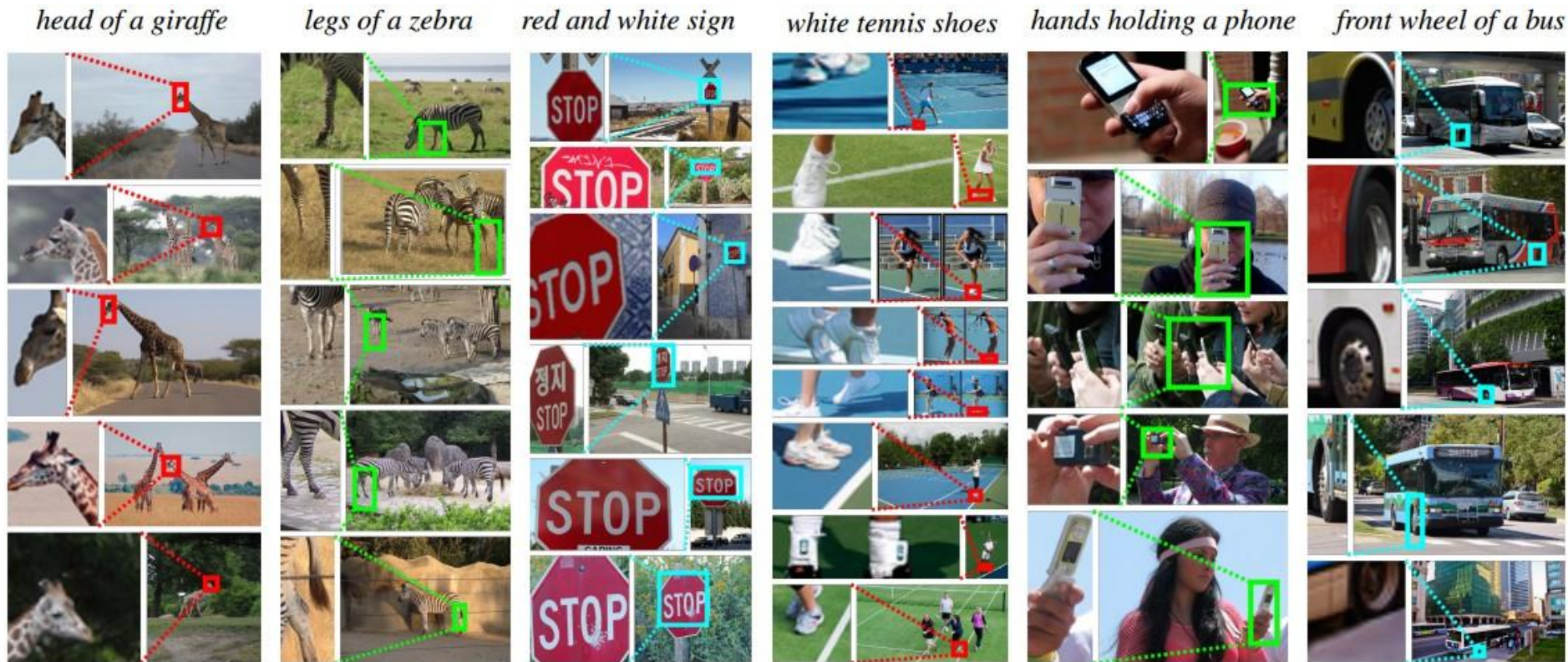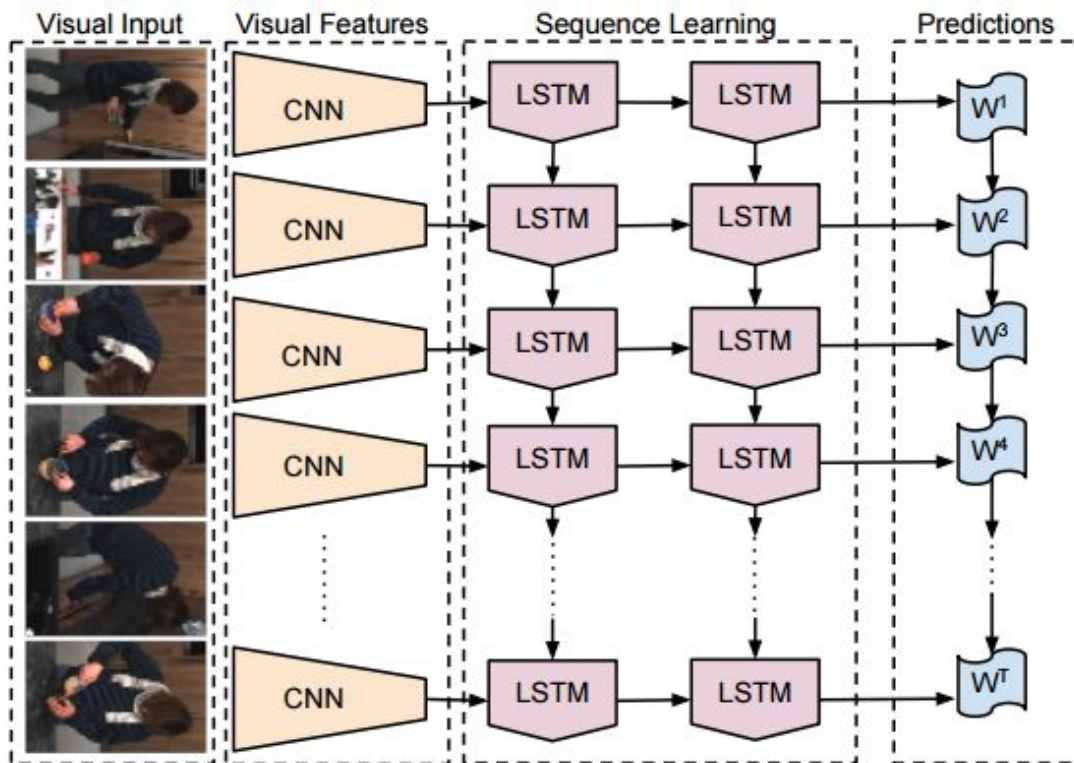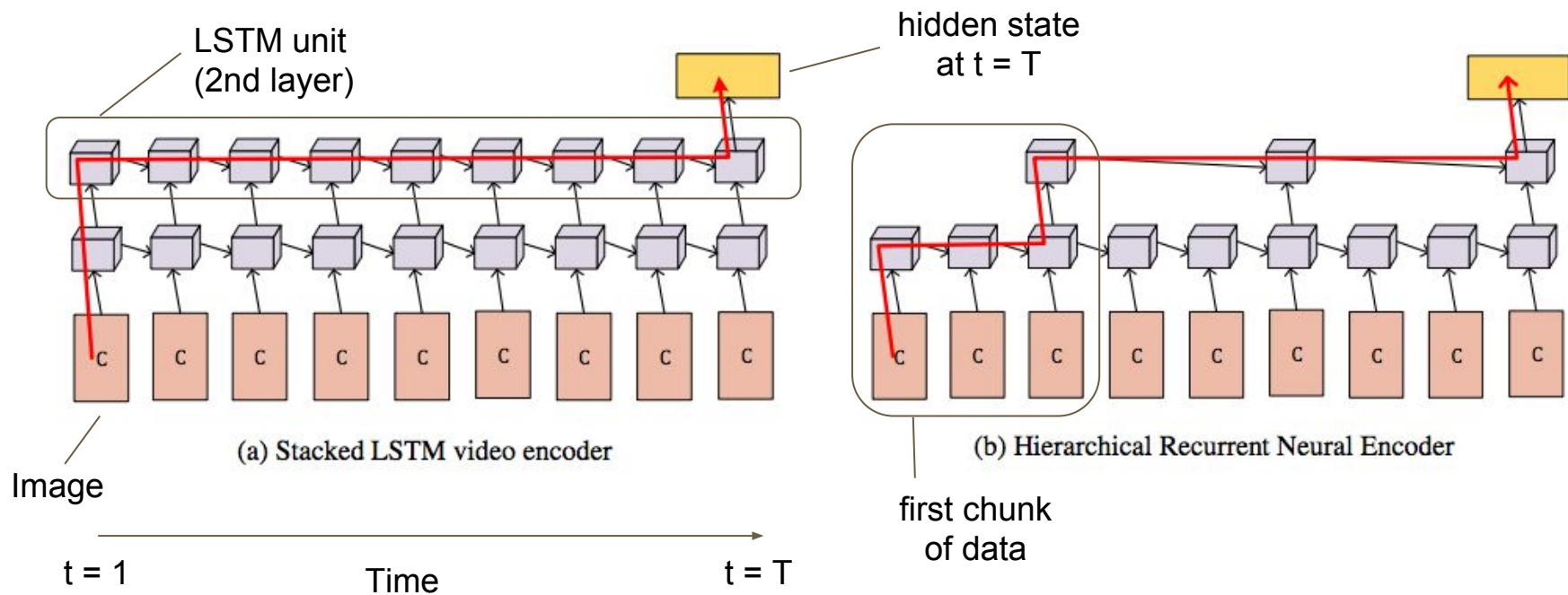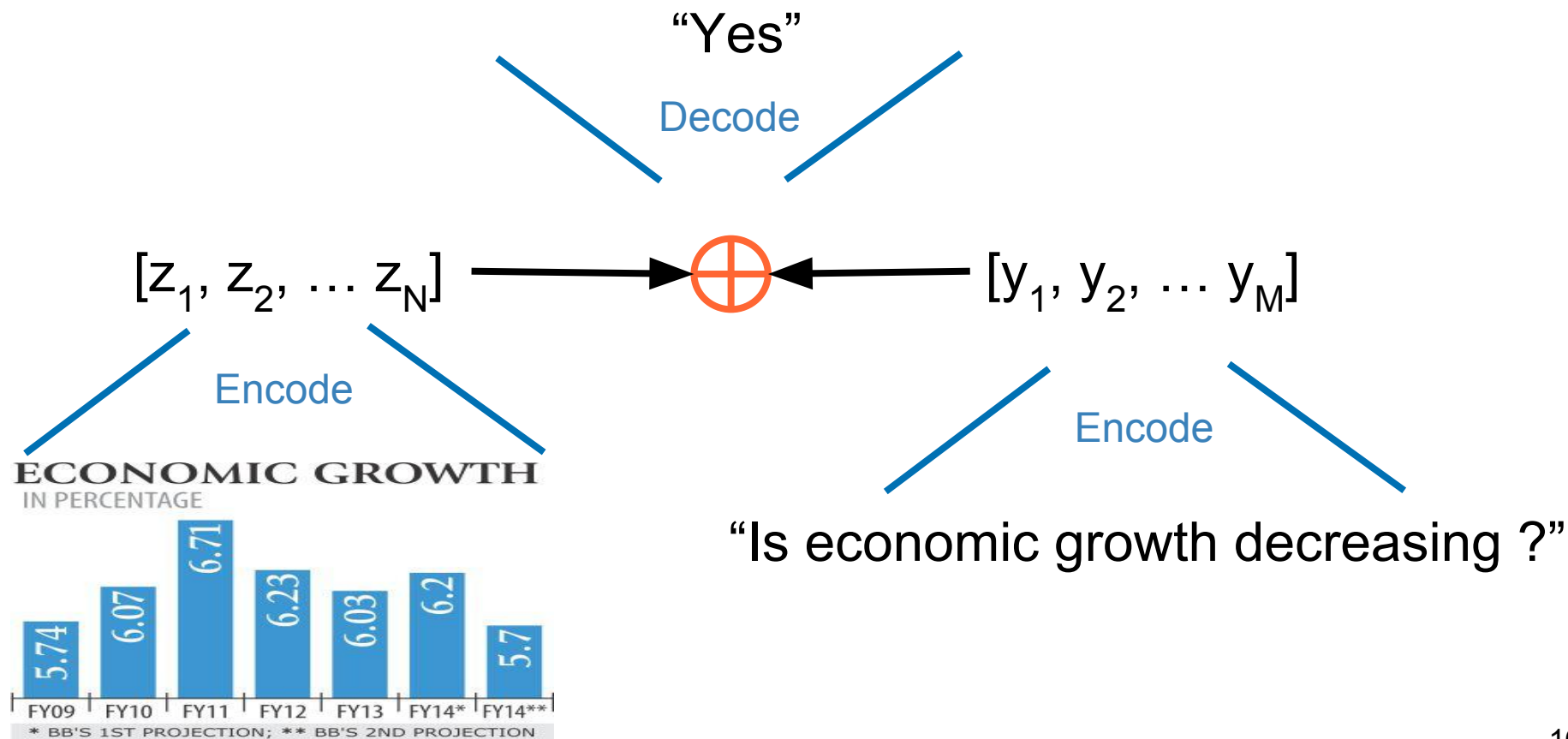
# Captioning: Video



Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, Trevor Darrel. Long-term Recurrent Convolutional Networks for Visual Recognition and Description, CVPR 2015. code

14

# Captioning: Video



(a) Stacked LSTM video encoder

(b) Hierarchical Recurrent Neural Encoder

LSTM unit (2nd layer)

hidden state at t = T

Image

first chunk of data

t = 1

Time

t = T

( Slides by Marc Bolaños) Pingbo Pan, Zhongwen Xu, Yi Yang,Fei Wu,Yueting Zhuang Hierarchical Recurrent Neural Encoder for Video Representation with Application to Captioning, CVPR 2016.

15

# Visual Question Answering

"Yes"

Decode

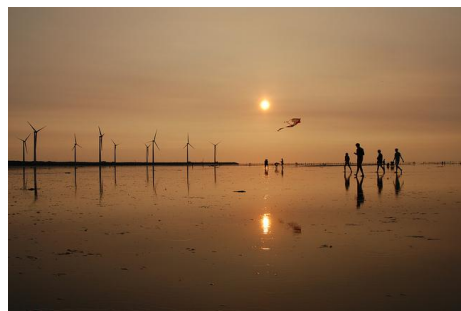$[z_1, z_2, \ldots z_N] \longrightarrow \oplus \longleftarrow [y_1, y_2, \ldots y_M]$

Encode



ECONOMIC GROWTH
IN PERCENTAGE

Encode

"Is economic growth decreasing ?"

# Visual Question Answering



What is the mustache made of?

bananas

Antol, Stanislaw, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. "VQA: Visual question answering." CVPR 2015.

17

# Visual Question Answering
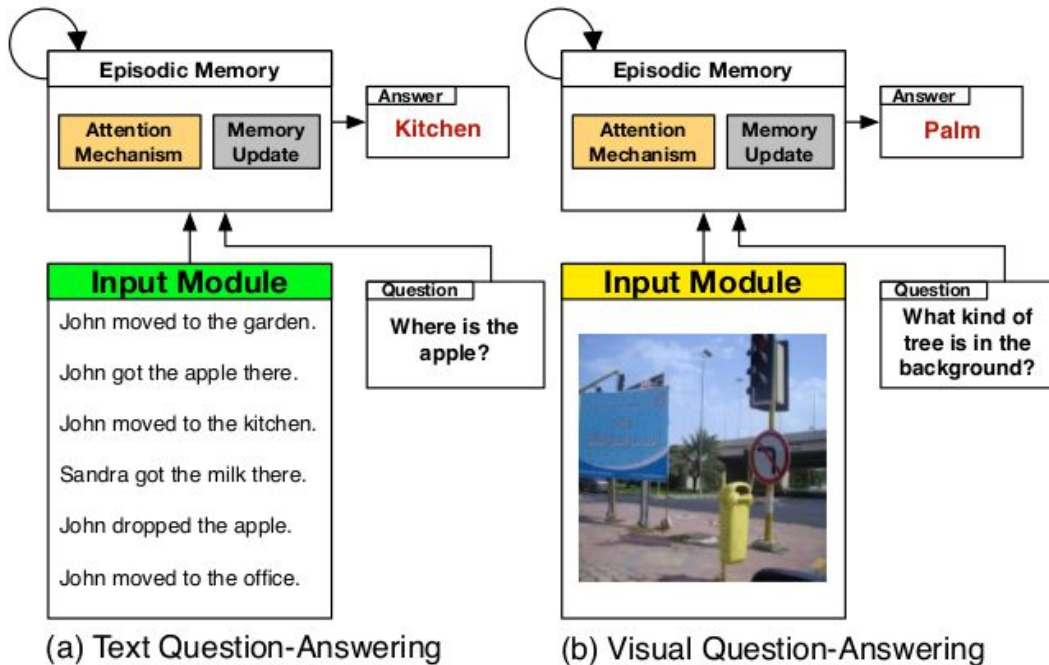


Slide credit: Issey Masuda

# Visual Question Answering



**Dynamic Parameter Prediction Network (DPPnet)**

Noh, H., Seo, P. H., & Han, B. Image question answering using convolutional neural network with dynamic parameter prediction. CVPR 2016

# Visual Question Answering: Dynamic



(a) Text Question-Answering

(b) Visual Question-Answering

# Visual Question Answering: Dynamic

**Main idea:** split image into local regions. Consider **each region equivalent to a sentence.**

**Local Region Feature Extraction:** CNN (VGG-19):
(1) Rescale input to 448x448.
(2) Take output from last pooling layer ➡ D=512x14x14 ➡ 196 512-d local region vectors.

**Visual feature embedding: W** matrix to project image features to *"q"-textual* space.



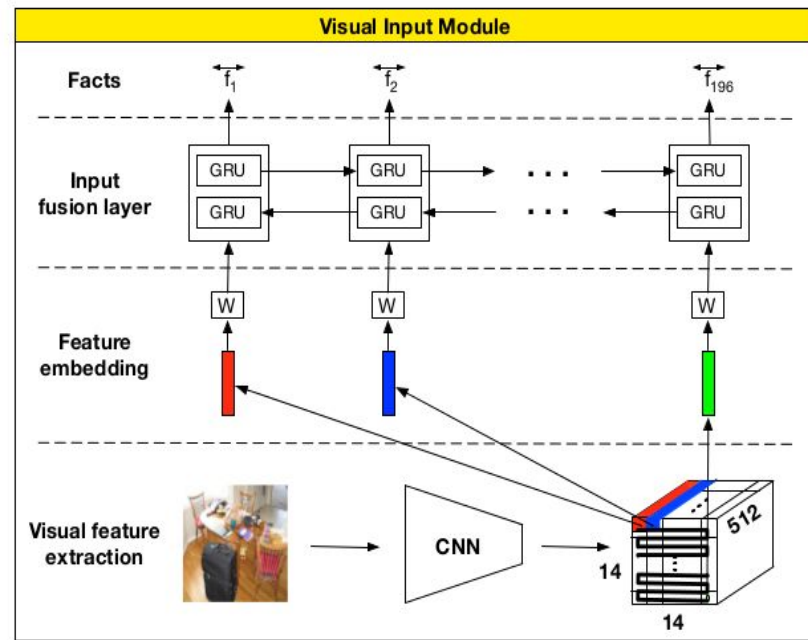*Figure 3.* VQA input module to represent images for the DMN.

21

# Visual Question Answering: Grounded

# Visual Dialog (Image Guessing Game)



Das, Abhishek, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra.
"Visual Dialog." CVPR 2017

23

# Visual Dialog (Image Guessing Game)



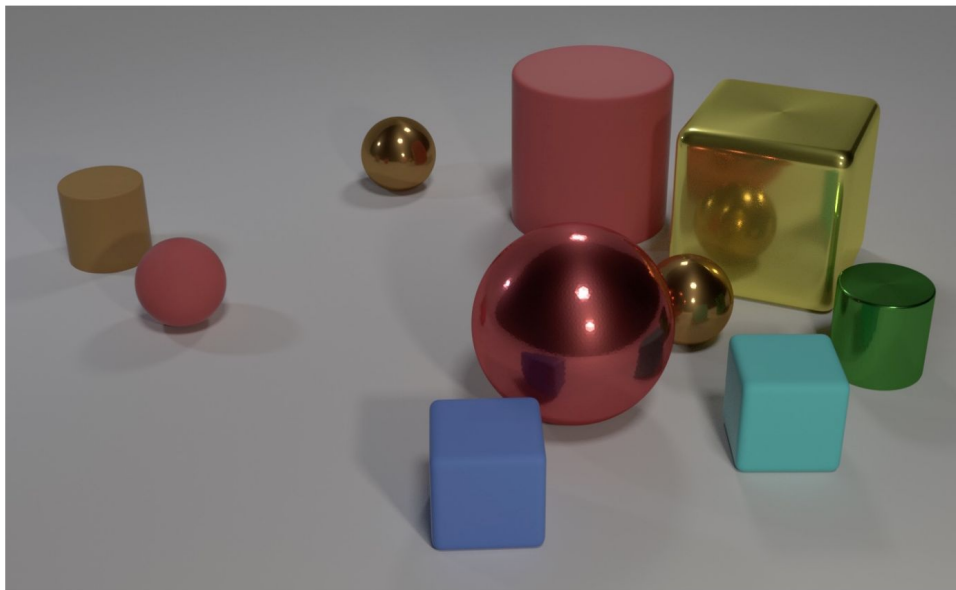Das, Abhishek, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra.
"Visual Dialog." CVPR 2017

24

# Visual Reasoning



**Q:** Are there an equal number of large things and metal spheres?
**Q:** What size is the cylinder that is left of the brown metal thing that is left of the big sphere? **Q:** There is a sphere with the same size as the metal cube; is it made of the same material as the small red sphere?
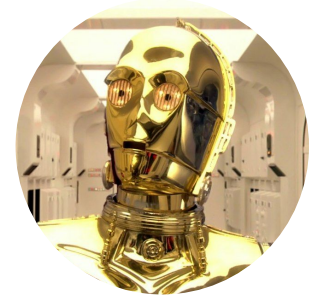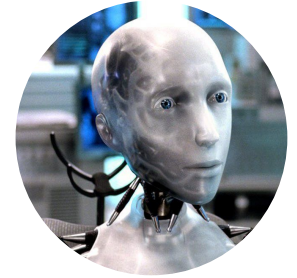
Johnson, Justin, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. "CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning." CVPR 2017

# Conclusions

New Turing test? How to evaluate AI's image understanding?

Slide credit: Issey Masuda
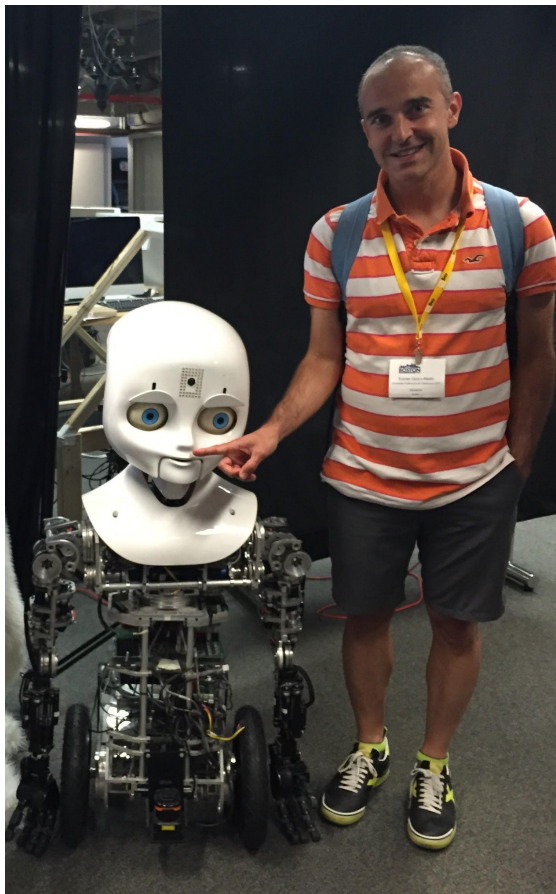
26

# Learn more

Julia Hockenmeirer



Flickr8k
[Hodosh et al. 2013]
~8,000 Flickr images
>40,000 captions
5 captions/image

# Thanks ! Q&A ?



Follow me at

 /ProfessorXavi

 @DocXavi

UNIVERSITAT POLITÈCNICA DE CATALUNYA
**BARCELONATECH**
**UPC**
**Department of Signal Theory and Communications**
*Image Processing Group*

https://imatge.upc.edu/web/people/xavier-giro