# Visual Question Answering with Deep Learning
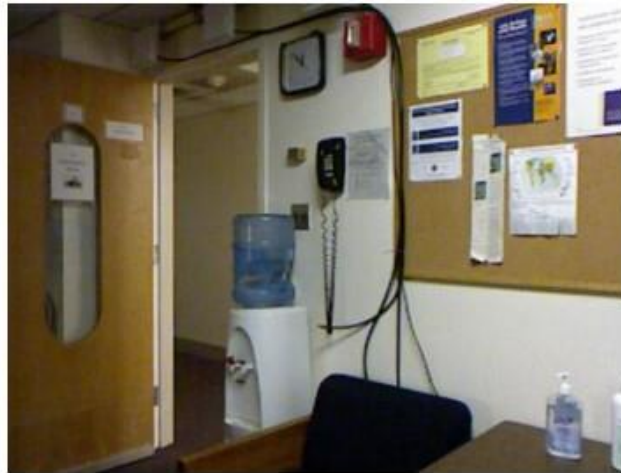
Aahitagni Mukherjee

Shubham Agarwal

# Visual Question Answering

Given an underline{image}, and a natural language-like underline{question}, find the correct underline{answer} to it

- Training on a set of triplets (image, question, answer).
- *Free-form* and *open-ended* questions.
- Answers can be single word or multiple word.



**Question:** what is the largest blue object in this picture?
**Ground truth:** water carboy
**Proposed CNN:** water carboy

**Question:** what color is the shade of the table lamp close to the bookshelf?
**Ground truth:** white
**Proposed CNN:** white

Lin Ma et al 2015

# Datasets

- DAQUAR(DAtaset for QUestion Answering on Real-world images) – 1450 images and 12468 questions related to them. On an average 12 words per question.

```
1   what is on the left side of the white oven on the floor and on right side of the blue armchair in the image1 ?
2   garbage_bin
3   what is on the left side of the fire extinguisher and on the right side of the chair in the image1 ?
4   table
5   what is between the the two white and black garbage bins in the image1 ?
6   chair
7   how many objects are between the fire extinguisher and the white oven on the floor in the image1 ?
8   3
9   what is the largest object in this picture in the image1 ?
10  washing_machine
```

- VQA(Visual Question Answering) dataset – 254,721 images, 764,163 questions, 9,934,119 answers

- Wu-Palmer Similarity Measure(WUPS score) is used for performance evaluation – Script by Malinowski M.
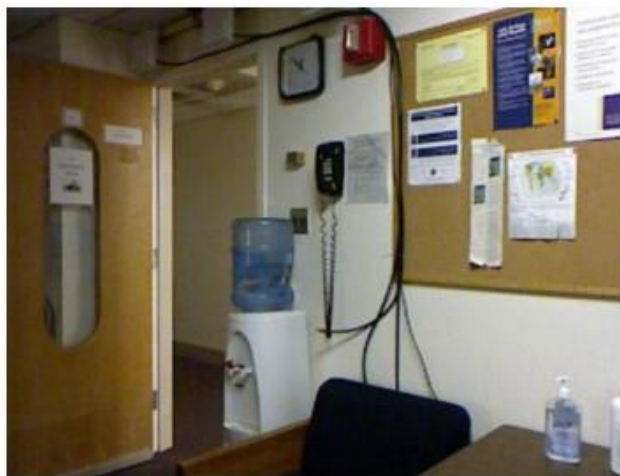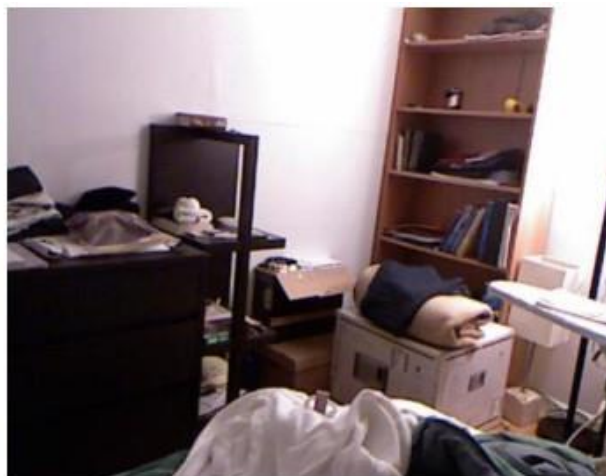
Malinowski et al. 2014
Anton et al. 2015
Wu et al. 1994

# Challenges

- The output is to be conditioned on both image and language inputs.
- A better representation of the image content is essential



**Question:** what is the largest blue object in this picture?
**Ground truth:** water carboy
**Proposed CNN:** water carboy

**Question:** what color is the shade of the table lamp close to the bookshelf?
**Ground truth:** white
**Proposed CNN:** white

- Interactions between the two modalities need to appropriately modelled.

# Previous Approaches

- **Neural-based approach** - image representation from a CNN is fed to each hidden layer of a single LSTM. The LSTM then models the concatenation of question and answer.

- **mQA approach** – 4 units - an LSTM to extract the *question representation*, a CNN to extract the *visual representation*, an LSTM for storing the linguistic context in an answer, and a *fusing component* to combine the information from the first three components and generate the answer.

- **VIS + LSTM** - Here the image is treated as a single word, and the intermediate representation of the input thus obtained is used for classification into the correct class, which is the single word answer.

- **CNN approach** - uses 3 CNN's - one to extract sentence representation, one for image representation, and the third is a multimodal layer to fuse the two.
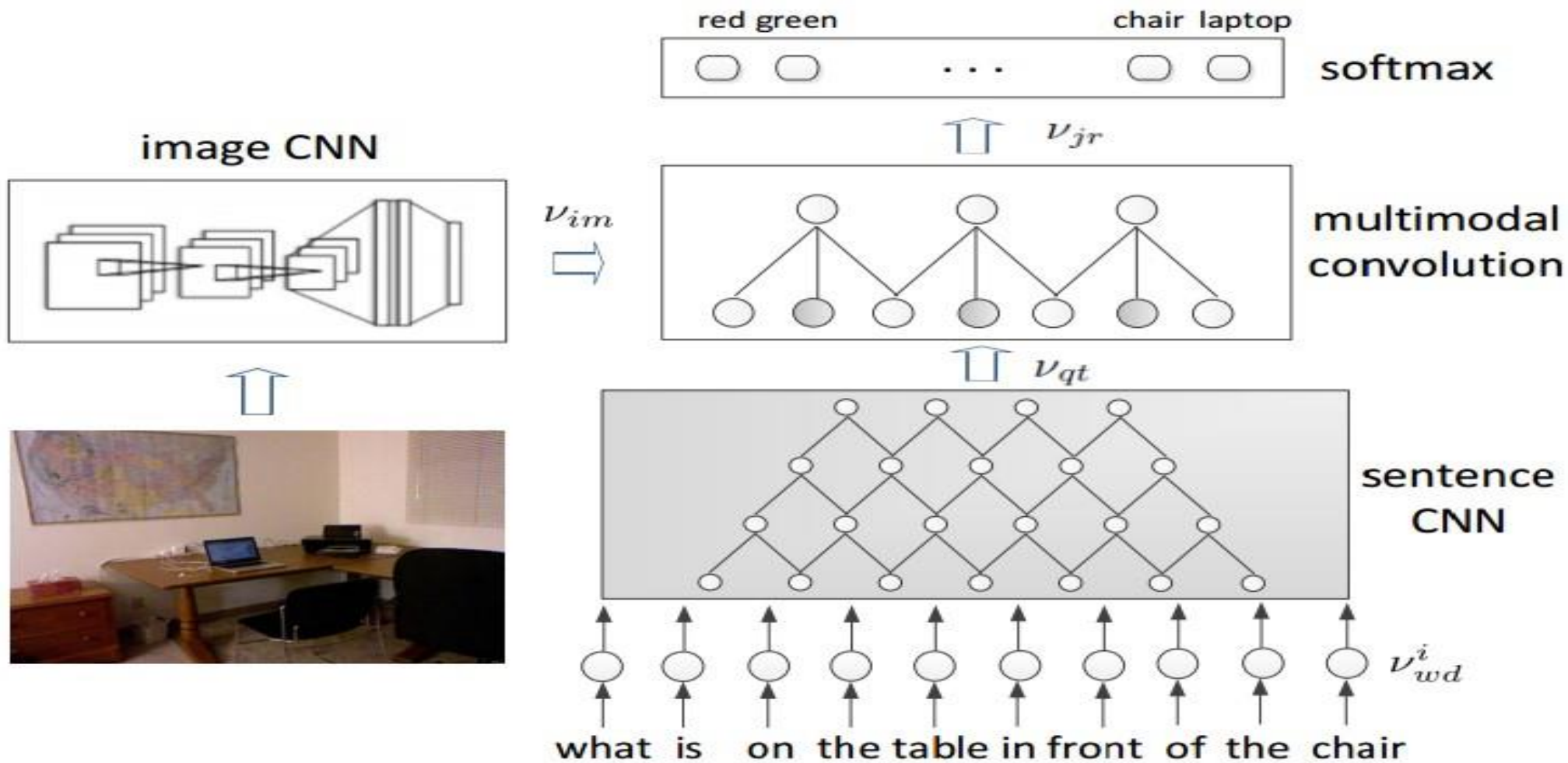
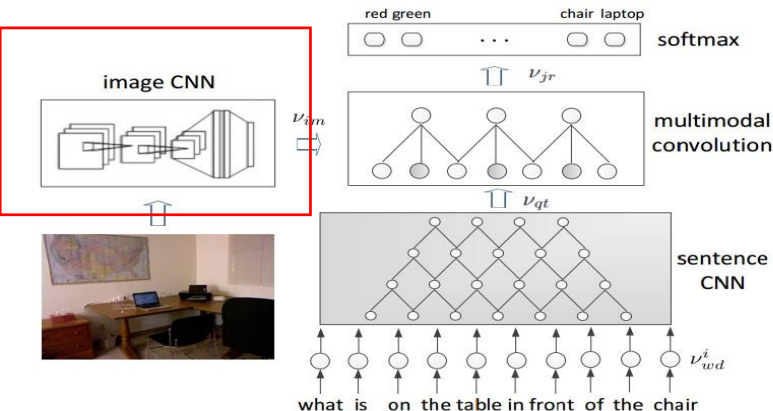Malinowski et al. 2015
Gao et al. 2015
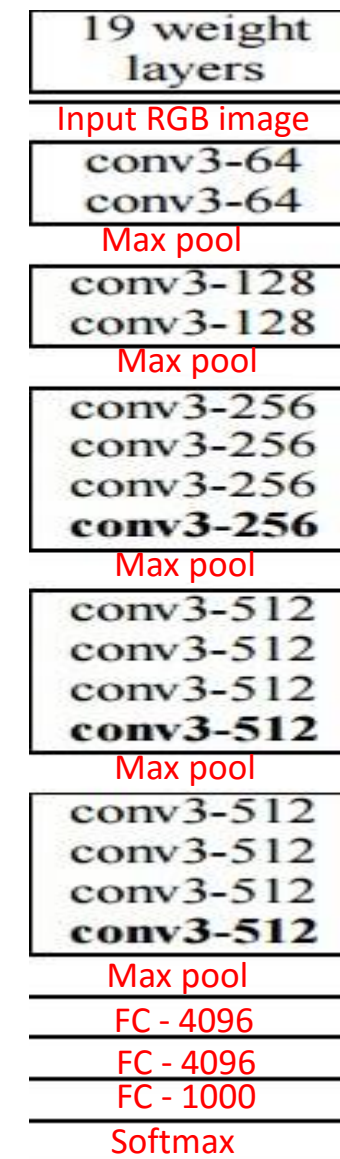Kiros et al. 2015
Lin Ma et al. 2015

# CNN model


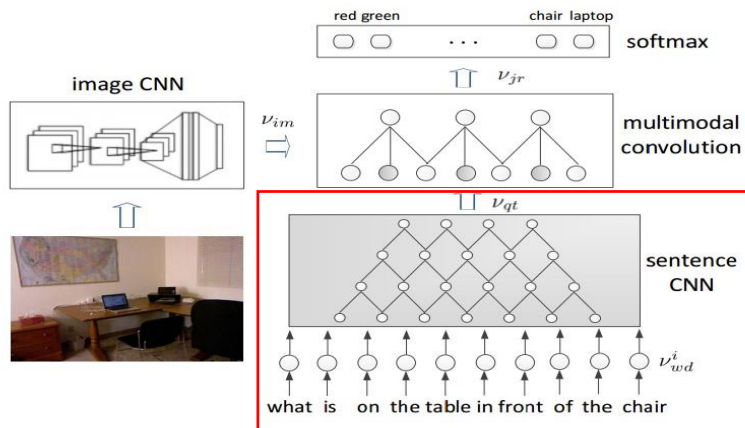
Lin Ma et al 2015

# Image CNN



$$\nu_{im} = \sigma(\boldsymbol{w_{im}}(CNN_{im}(I)) + b_{im})$$

$\sigma$: Nonlinear activation function.

$\boldsymbol{w_{im}}|_{d \times 4096}$ : Mapping matrix

$CNN_{im}$ takes image as input and outputs a fixed length vector.

| 19 weight layers |
|---|
| Input RGB image |
| conv3-64 |
| conv3-64 |
| Max pool |
| conv3-128 |
| conv3-128 |
| Max pool |
| conv3-256 |
| conv3-256 |
| conv3-256 |
| **conv3-256** |
| Max pool |
| conv3-512 |
| conv3-512 |
| conv3-512 |
| **conv3-512** |
| Max pool |
| conv3-512 |
| conv3-512 |
| conv3-512 |
| **conv3-512** |
| Max pool |
| FC - 4096 |
| FC - 4096 |
| FC - 1000 |
| Softmax |

Simoyan et al. 2015

# Sentence CNN

image CNN

red green    chair laptop

softmax

$\nu_{jr}$

$\nu_{im}$

multimodal convolution

$\nu_{qt}$

sentence CNN

$\nu_{wd}^i$

what is on the table in front of the chair

$\nu_{(\ell)}^i$

$\nu_{(\ell-1)}^i$    $\nu_{(\ell-1)}^{i+1}$    $\nu_{(\ell-1)}^{i+2}$

**1** For sequential input $\sigma$, convolution unit for feature map of type $f$ on the $l^{th}$ layer is

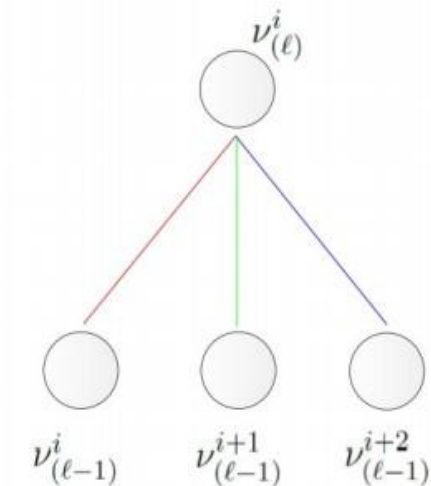$$\nu_{(l,f)}^i \overset{\text{def}}{=} \sigma(\boldsymbol{w}_{(l,f)} \vec{\nu}_{(l-1)}^i + b_{(l,f)})$$

**2** $\vec{\nu}_{(l-1)}^i \overset{\text{def}}{=} \nu_{(l-1)}^i \mid\mid \nu_{(l-1)}^{i+1} \mid\mid \nu_{l-1}^{i+2}$

**3** $\vec{\nu}_{(0)}^i \overset{\text{def}}{=} \nu_{wd}^i \mid\mid \nu_{wd}^{i+1} \mid\mid \nu_{wd}^{i+2}$

**4** Max-pooling after each convolution

$$\nu_{(l+1,f)}^i = max(\nu_{(l,f)}^{2i}, \nu_{(l,f)}^{2i+1})$$

$\nu_{wd}^i$ : Skip-gram    word embedding of i-th question word

Lin Ma et al 2015

# Multimodal Convolutional Layer



red green        chair laptop

$\dots$        softmax

image CNN

multimodal convolution
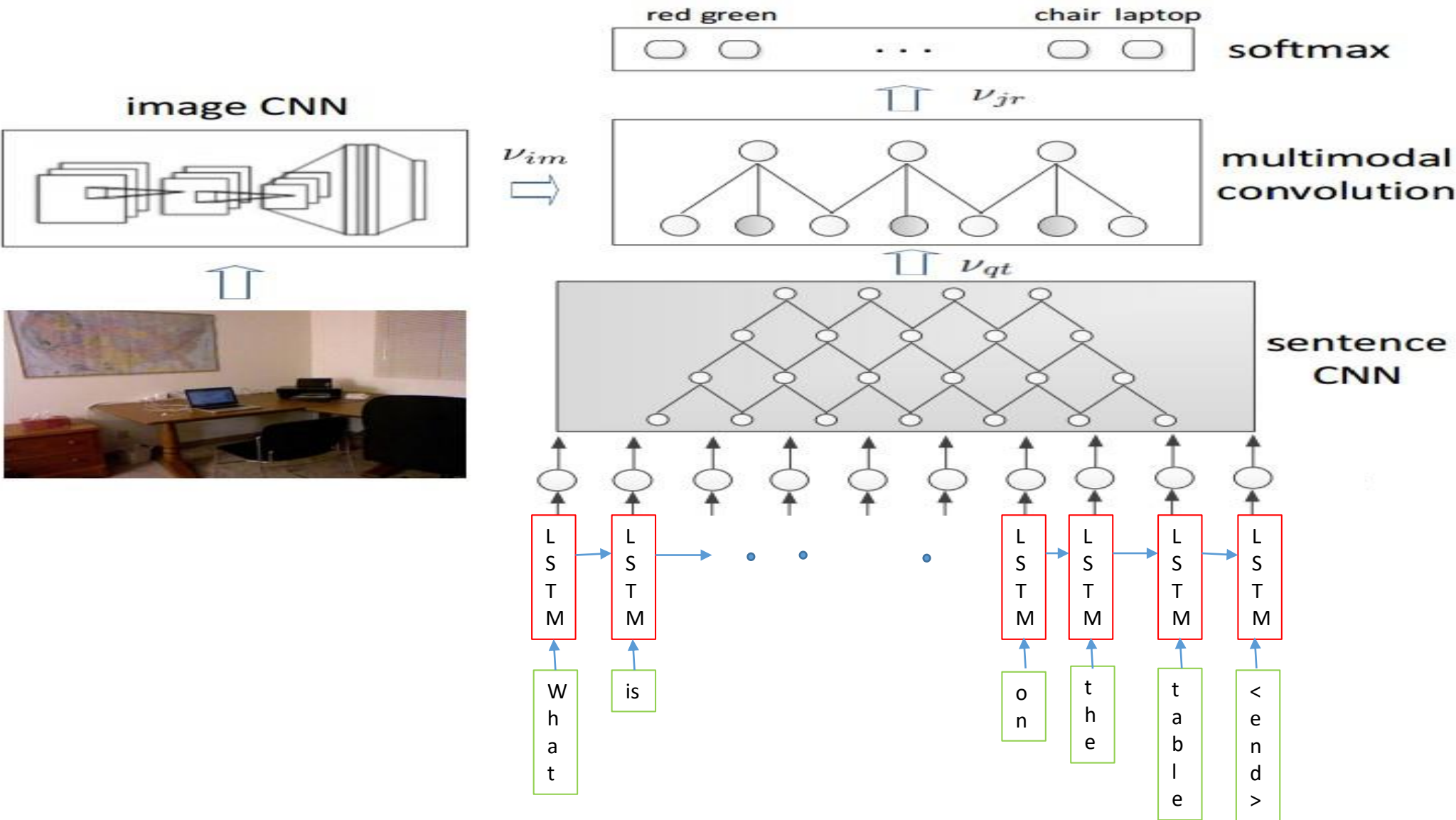
sentence CNN

what is on the table in front of the chair

input: $\nu_{qt} = [\nu_{(6)}^0 \dots \nu_{(6)}^n]$

Capturing the interaction between two multimodal inputs

$$\vec{\nu}_6^i = \nu_6^i \,||\, \nu_{im} \,||\, \nu_6^{i+1}$$

$$\nu_{(mm,f)}^i = \sigma(\boldsymbol{w}_{(mm,f)}\vec{\nu}_{(6)}^i + b_{(mm,f)})$$

Lin Ma et al 2015

# Proposed Modification

# Input to LSTM

Skip-gram word embeddings from the question sentence

Mikolov et al. 2013
https://code.google.com/p/word2vec/

# LSTM

## LSTM Unit



$$i_t = \sigma(W_{vi}v_t + W_{hi}h_{t-1} + b_i)$$
$$f_t = \sigma(W_{vf}v_t + W_{hf}h_{t-1} + b_f)$$
$$o_t = \sigma(W_{vo}v_t + W_{ho}h_{t-1} + b_o)$$
$$g_t = \phi(W_{vg}v_t + W_{hg}h_{t-1} + b_g)$$
$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \phi(c_t)$$



Malinowski et al. 2015
https://github.com/junhyukoh/caffe-lstm

# References

1. [Lin Ma et al. 2015] L. Ma and Z. Lu and H. Li. Learning to Answer Questions From Image using Convolutional Neural Network 2015.  arXiv:1506.00333.
2. [Malinowski et al. 2014] M. Malinowski and M. Fritz 2014b. Towards a Visual Turing Challenge. arXiv:1410.8027.
3. [Anton et al. 2015] S. Antol and A. Agrawal and J. Lu and M. Mitchell and D. Batra and C. Lawrence Zitnick and D. Parikh 2015. VQA: Visual Question Answering. arXiv:1505.00468.
4. [Wu et al. 1994] Z. Wu and M. Palmer 1994. Verb Semantics and Lexical Selection. ACL 1994.
5. [Gao et al. 2015] H. Gao and J. Mao and J. Zhou and Z. Huang and L. Wang and W. Xu 2015. Are You Talking to a Machine? Dataset and Methods for Multilingual Image Question Answering. arXiv:1505.05612.
6. [Kiros et al. 2015] M Ren, R. Kiros, R. Zemel. Exploring Models and Data for Image Question Answering 2015. arXiv:1505.02074.
7. [Simonyan et al. 2015] K. Simonyan and A. Zisserman 2014. Very Deep Convolutional Networks for Large-scale Image Recognition. arXiv 1409.1556.
8. [Mikolov et al. 2013] T. Mikolov and K. Chen and G. Corrado and J. Dean 2013. Efficient Estimation of Word Representations in Vector Space. arXiv:1301.3781.
9. [Malinowski et al. 2015] M. Malinowski and M. Rohrbach and M. Fritz 2015b. Ask Your Neurons: A Neural-based Approach to Answering Questions about Images. arXiv 1505.01121.

# CONVOLUTIONAL, LONG SHORT-TERM MEMORY, FULLY CONNECTED DEEP NEURAL NETWORKS

*Tara N. Sainath, Oriol Vinyals, Andrew Senior, Haşim Sak*

Google, Inc., New York, NY, USA
{tsainath, vinyals, andrewsenior, hasim}@google.com

## ABSTRACT

Both Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) have shown improvements over Deep Neural Networks (DNNs) across a wide variety of speech recognition tasks. CNNs, LSTMs and DNNs are complementary in their modeling capabilities, as CNNs are good at reducing frequency variations, LSTMs are good at temporal modeling, and DNNs are appropriate for mapping features to a more separable space. In this paper, we take advantage of the complementarity of CNNs, LSTMs and DNNs by combining them into one unified architecture. We explore the proposed architecture, which we call CLDNN, on a variety of large vocabulary tasks, varying from 200 to 2,000 hours. We find that the CLDNN provides a 4-6% relative improvement in WER over an LSTM, the strongest of the three individual models.

nonlinear hidden layer. If factors of variation in the hidden states could be reduced, then the hidden state of the model could summarize the history of previous inputs more efficiently. In turn, this could make the output easier to predict. Reducing variation in the hidden states can be modeled by having DNN layers after the LSTM layers. This is similar in spirit to the hidden to output model proposed in [4], and also tested for speech, though with RNNs [8].

The model we propose is to feed input features, surrounded by temporal context, into a few CNN layers to reduce spectral variation. The output of the CNN layer is then fed into a few LSTM layers to reduce temporal variations. Then, the output of the last LSTM layer is fed to a few fully connected DNN layers, which transform the features into a space that makes that output easier to classify.

Combining CNN, LSTMs and DNNs has been explored in [9]. However, in that paper the three models were first trained separately and then the three outputs were combined through a combination

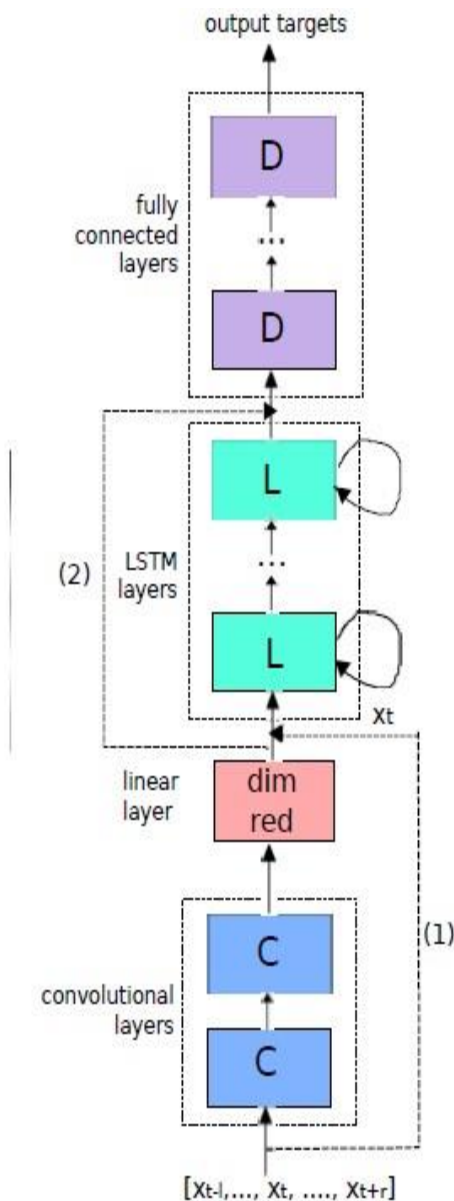| Method | WER |
|----------|------|
| LSTM | 18.0 |
| CNN+LSTM | 17.6 |
| LSTM+DNN | 17.6 |
| CLDNN | **17.3** |

**Table 5.** WER, CLDNN



**Fig. 1.** CLDNN Architecture