

# MULTILINGUAL IMAGE DESCRIPTION WITH NEURAL SEQUENCE MODELS

---

Desmond Elliott, Stella Frank, Eva Hasler

February 3, 2016

iV&L Net Working Group 3 Meeting



INSTITUTE FOR LOGIC, LANGUAGE AND COMPUTATION

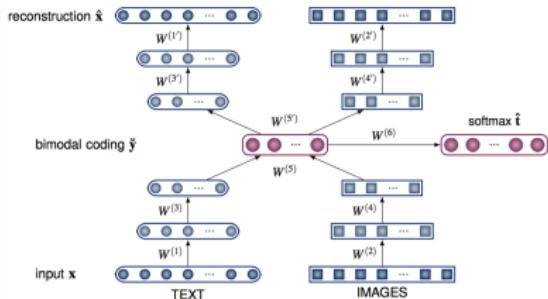


Stella Frank

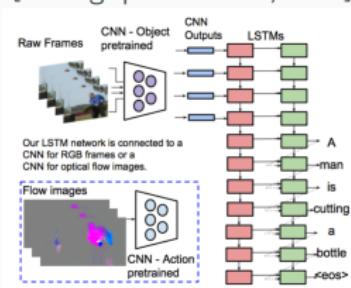


Eva Hasler

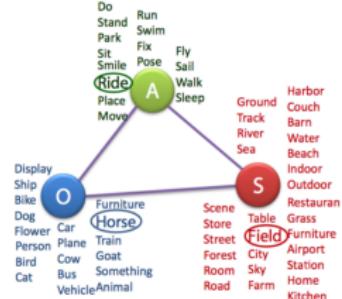
## Grounded Semantics [Silberer and Lapata, 2014]



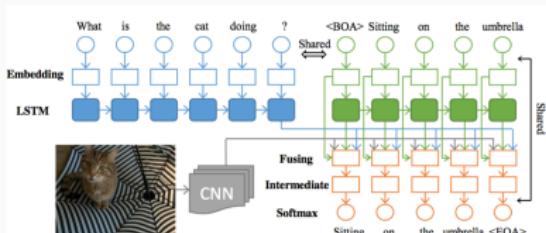
## Video Description [Venugopalan et al., 2015]



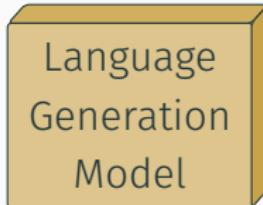
## Image Description [Farhadi et al., 2010]



## Question-Answering [Gao et al., 2015]



## BRIEFEST OVERVIEW OF IMAGE DESCRIPTION

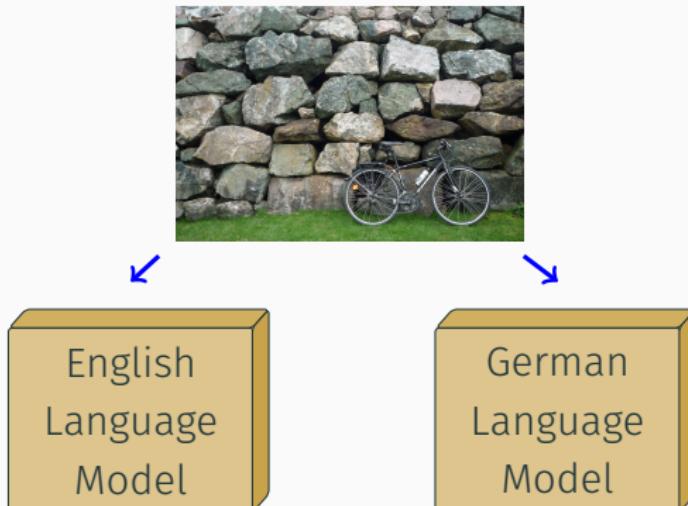


A bike is leaning against a stone wall

See Bernardi et al. [2016] for an overview of datasets, models, and evaluations.

## THIS TALK: DESCRIBING IMAGES IN MULTIPLE LANGUAGES

- Extend image description generation to new languages
- Text-based image search in any language
- Localised alt-text generation on the Web
- Translate movie descriptions



## HOW CAN WE EXPLOIT MULTILINGUAL MULTIMODAL CONTEXT?

Ein Rad steht neben der Mauer → A bicycle / wheel ...



## HOW CAN WE EXPLOIT MULTILINGUAL MULTIMODAL CONTEXT?

Ein **Rad** steht neben der Mauer → A bicycle / wheel ...

<Image features> → A ? is leaning against the wall



## HOW CAN WE EXPLOIT MULTILINGUAL MULTIMODAL CONTEXT?

Ein Rad steht neben der Mauer → A bicycle / wheel ...

<Image features> → A ? is leaning against the wall

Possible solutions:

- Collect more data



# HOW CAN WE EXPLOIT MULTILINGUAL MULTIMODAL CONTEXT?

Ein Rad steht neben der Mauer → A bicycle / wheel ...

<Image features> → A ? is leaning against the wall

Possible solutions:

- Collect more data
- Exploit data in a different modality (images or video)



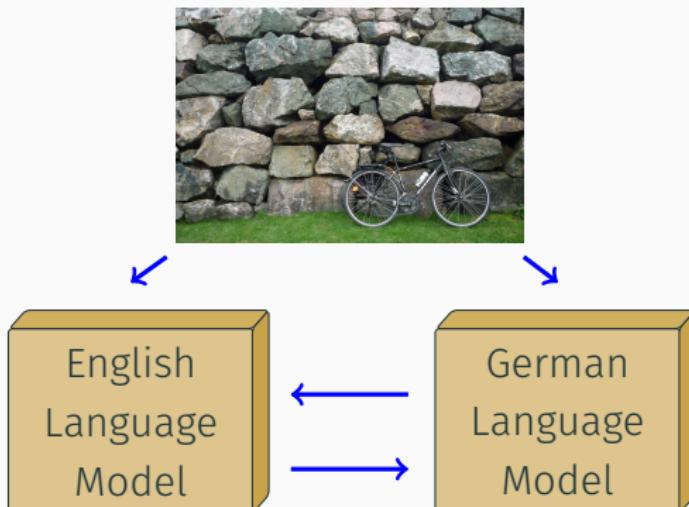
## HOW CAN WE EXPLOIT MULTILINGUAL MULTIMODAL CONTEXT?

Ein Rad steht neben der Mauer → A bicycle / wheel ...

<Image features> → A ? is leaning against the wall

Possible solutions:

- Collect more data
- Exploit data in a different modality (images or video)



Let  $t$  be the **target** language description,  $s$  be the **source** language description and  $i$  be the **image**.

### 1. Multimodal Machine Translation

- Always given a source description and image
- $\hat{t} = \operatorname{argmax}_t p(t|i, s)$

## TWO TASKS FOR MULTILINGUAL IMAGE DESCRIPTION

Let  $t$  be the **target** language description,  $s$  be the **source** language description and  $i$  be the **image**.

### 1. Multimodal Machine Translation

- Always given a source description and image
- $\hat{t} = \operatorname{argmax}_t p(t|i, s)$

### 2. Crosslingual Image Description

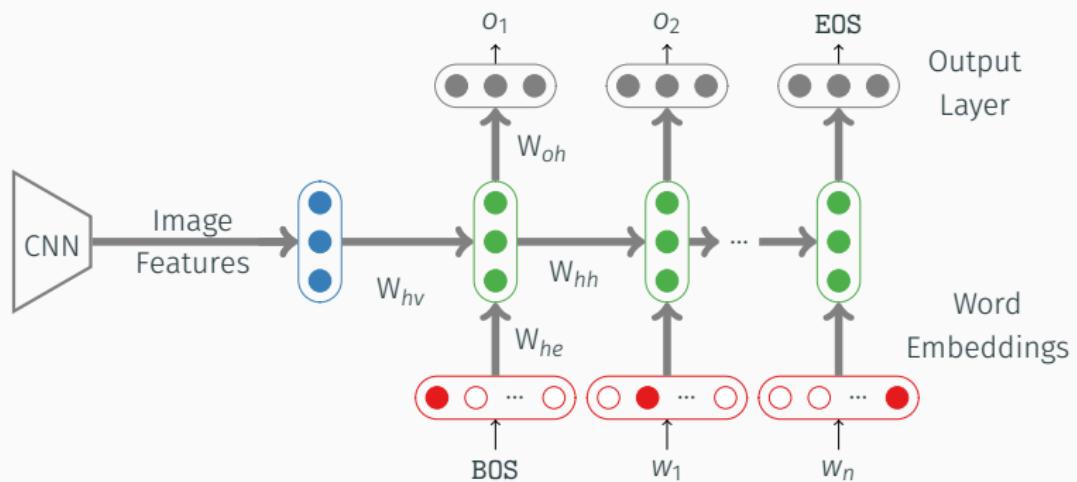
- Automatically generate a source language description
- $\hat{t} = \operatorname{argmax}_t p(t|i, \hat{s})$

## MULTILINGUAL MULTIMODAL MODEL

---

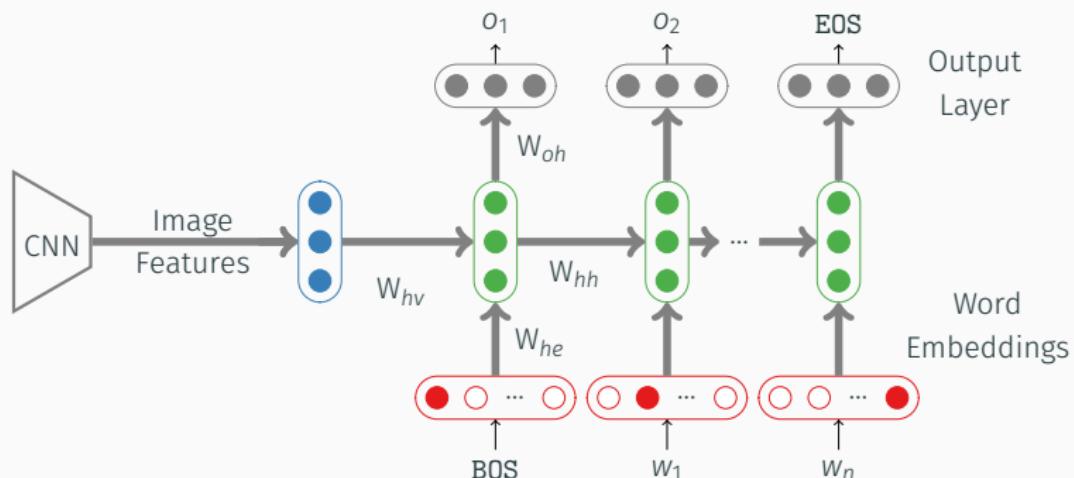
# MULTIMODAL LANGUAGE MODELS

[VINYALS ET AL., 2015, KARPATHY AND FEI-FEI, 2015]



# MULTIMODAL LANGUAGE MODELS

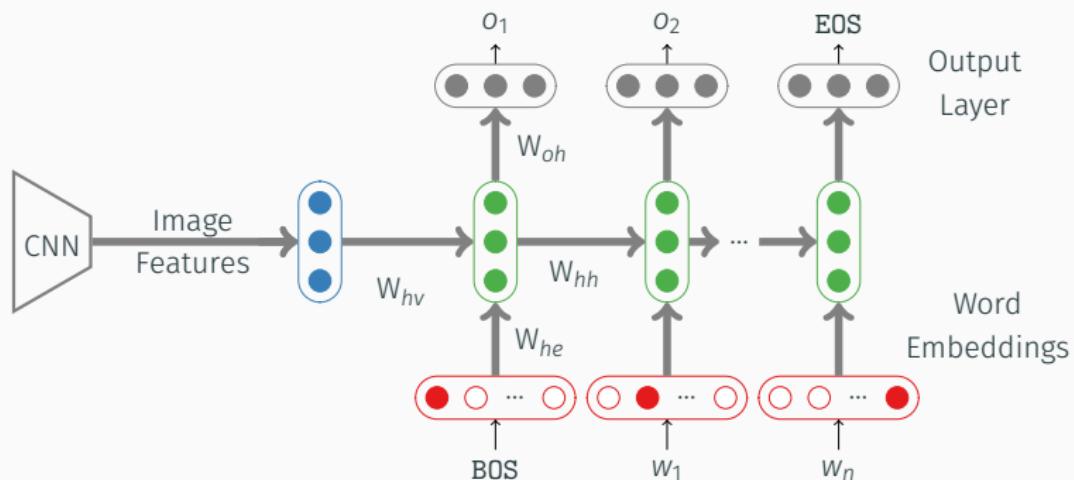
[VINYALS ET AL., 2015, KARPATHY AND FEI-FEI, 2015]



$$\cdot e_i = W_{he}w_i$$

# MULTIMODAL LANGUAGE MODELS

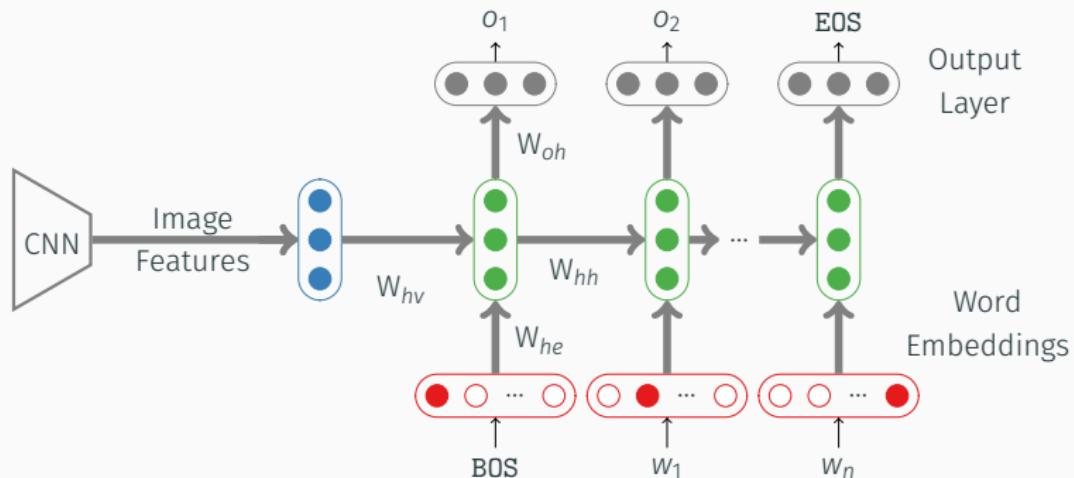
[VINYALS ET AL., 2015, KARPATHY AND FEI-FEI, 2015]



- $e_i = W_{he}w_i$
- $h_i = f(W_{hh}h_{i-1} + e_i + \mathbb{1}(t=0) \cdot W_{hv}v)$

# MULTIMODAL LANGUAGE MODELS

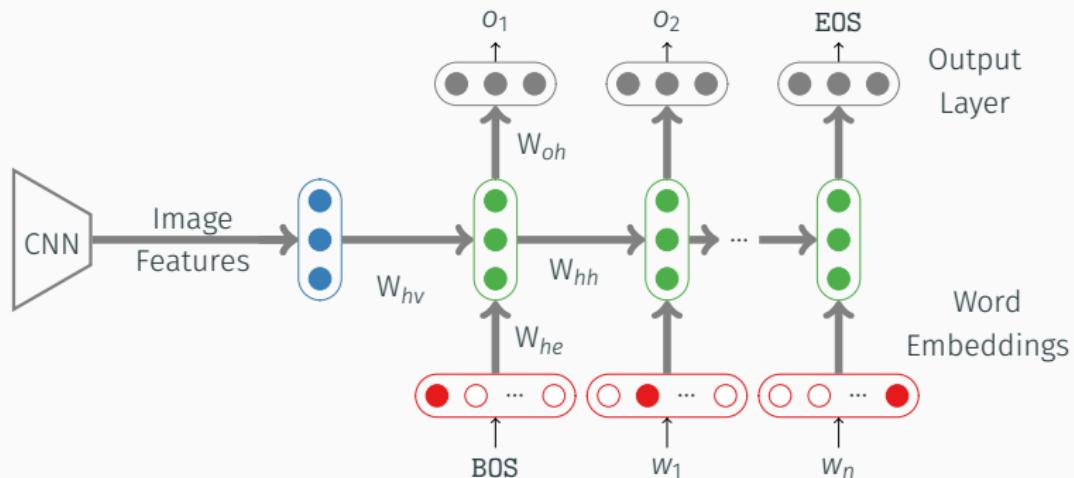
[VINYALS ET AL., 2015, KARPATHY AND FEI-FEI, 2015]



- $e_i = W_{he}w_i$
- $h_i = f(W_{hh}h_{i-1} + e_i + \mathbb{1}(t=0) \cdot W_{hv}v)$
- $o_i = \text{softmax}(W_{oh}h_i)$

# MULTIMODAL LANGUAGE MODELS

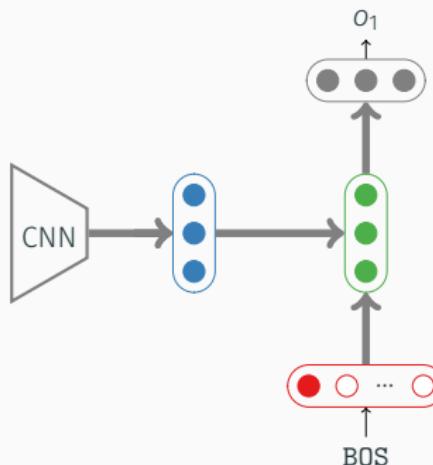
[VINYALS ET AL., 2015, KARPATHY AND FEI-FEI, 2015]



- $e_i = W_{he}w_i$
- $h_i = f(W_{hh}h_{i-1} + e_i + \mathbb{1}(t=0) \cdot W_{hv}v)$
- $o_i = \text{softmax}(W_{oh}h_i)$
- $\text{Loss} = - \sum_{n=1}^N \sum_{i=1}^K \log p(o_i)$

# INFERENCE WITH MULTIMODAL LANGUAGE MODELS

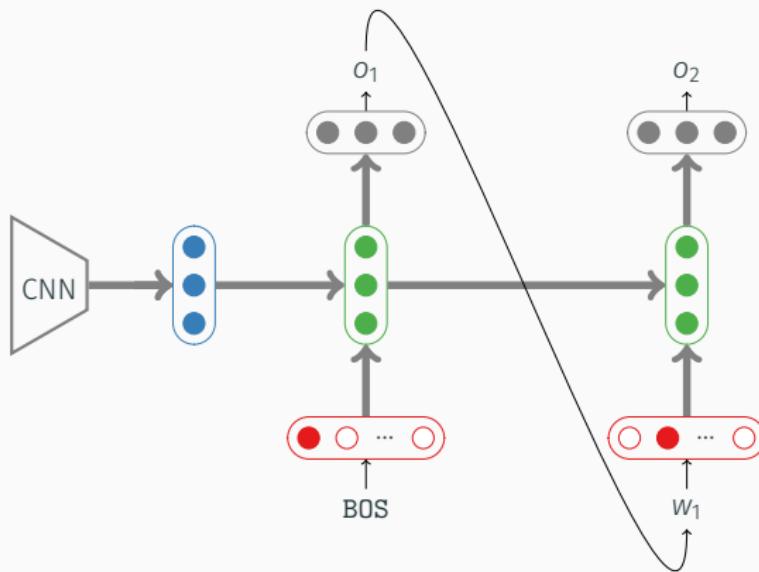
[VINYALS ET AL., 2015, KARPATHY AND FEI-FEI, 2015]



- Initialise with image features and BOS token

# INFERENCE WITH MULTIMODAL LANGUAGE MODELS

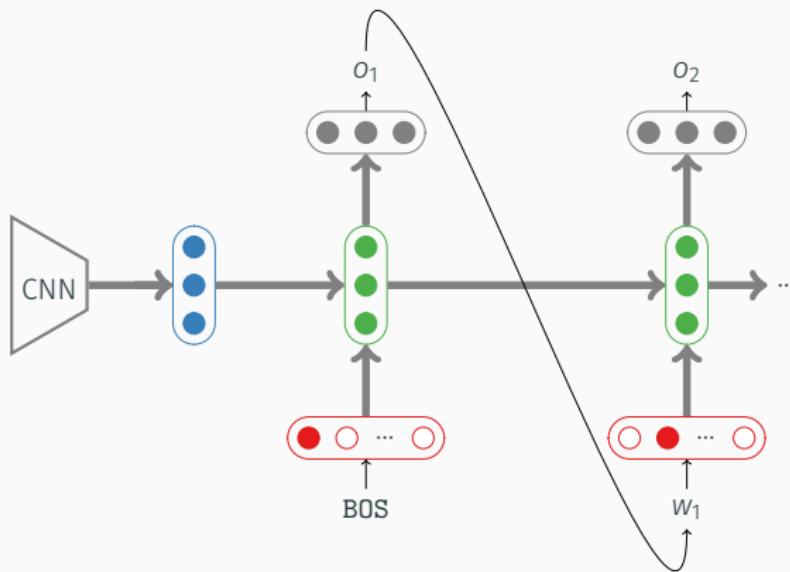
[VINYALS ET AL., 2015, KARPATHY AND FEI-FEI, 2015]



- Initialise with image features and BOS token
- Feed sampled word into the next timestep

# INFERENCE WITH MULTIMODAL LANGUAGE MODELS

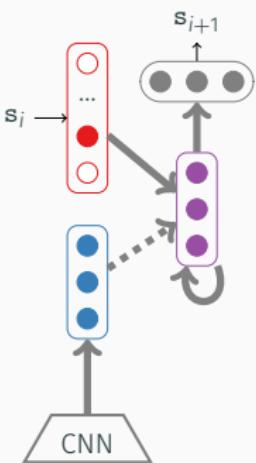
[VINYALS ET AL., 2015, KARPATHY AND FEI-FEI, 2015]



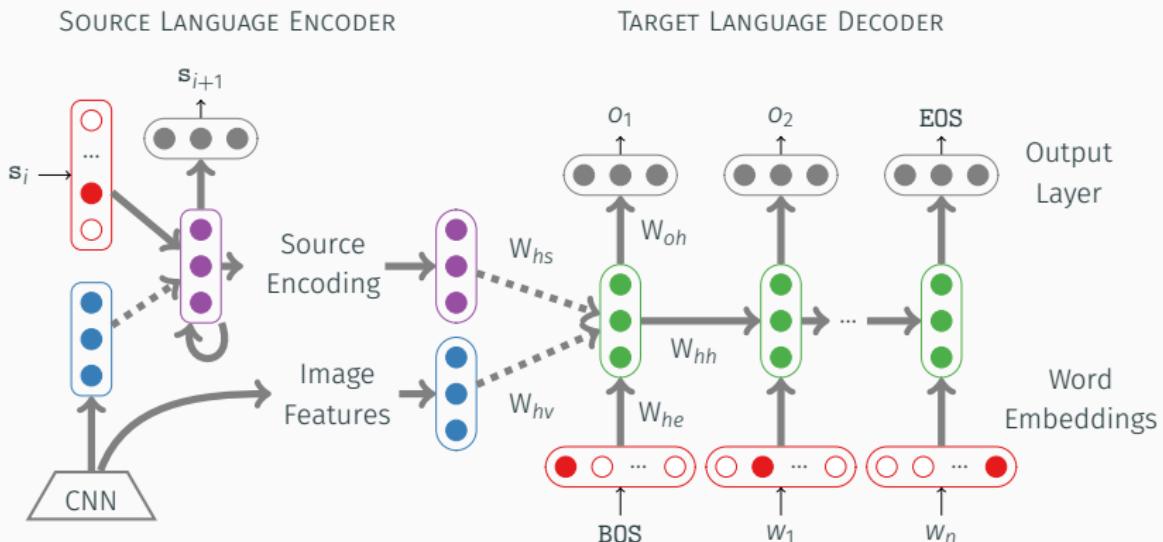
- Initialise with image features and BOS token
- Feed sampled word into the next timestep
- Decode until emit EOS token

# MULTILINGUAL MULTIMODAL MODEL [ELLIOTT ET AL., 2015]

SOURCE LANGUAGE ENCODER



# MULTILINGUAL MULTIMODAL MODEL [ELLIOTT ET AL., 2015]



$$h_i = f(W_{hh}h_{-1} + e_i + \mathbb{1}(t=0) \cdot W_{hv}v + \mathbb{1}(t=0) \cdot W_{hs}s)$$

## MULTILINGUAL MULTIMODAL MODEL (CONT.)

- Each model trained towards its own objective, unlike Sequence-to-Sequence Learning [Sutskever et al., 2014]
  - CNN: object recognition
  - Source LM: source language generation
  - Target LM: target language generation
- MMLM learns task-specific representations given transferred inputs
  - e.g. Target-LM with **multimodal** source features vs. separate **visual and source** features
  - Easily work on new languages with fixed input representations

## EXPERIMENTS

---

- Generate description in target language
  - Measures<sup>1</sup>: Meteor, BLEU, Perplexity
1. Multimodal Machine Translation
    - Always given a source description and image
  2. Crosslingual Image Description
    - Given an image, automatically generate a source description with a source MLM
      - and pass encoded textual features to a target LM
      - and pass encoded visual+textual features to a target LM

---

<sup>1</sup>See Elliott and Keller [2014] and Vedantam et al. [2015] for more details on measuring image description quality

1. a yellow building with white columns in the background
2. two palm trees in front of the house



1. ein gelbes Gebäude mit weißen Säulen im Hintergrund
2. zwei Palmen vor dem Haus

- 17,655 training / 1,962 testing
- Up to five semantically diverse descriptions / image
  - We use **only** the first description
- Descriptions translated from English to German

- Models are built using Keras library
- Adam optimiser [Kingma and Ba, 2014]
- Mini-batches of 100 examples
- Dropout over word, visual, and source features ( $p = 0.5$ )
- LSTM with 256-D memory cell [Hochreiter and Schmidhuber, 1997]
- 4096-D visual features from 15th layer of VGG-16 CNN [Simonyan and Zisserman, 2015]
- 256-D source feature vectors
- 256-D word embedding features
- Vocabulary size German: 2,374, English: 1,763 (UNK<3)

## MODELS: MULTIMODAL LANGUAGE MODEL (MLM)



Schulkinder sitzen  
in einem Klassenzimmer



TARGET MODEL

## MODELS: SOURCE LM → TARGET LM



children sitting  
in a classroom

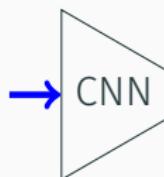
SOURCE MODEL

Schulkinder sitzen  
in einem Klassenzimmer



TARGET MODEL

## MODELS: SOURCE MLM → TARGET LM



children sitting  
in a classroom

SOURCE MODEL

Schulkinder sitzen  
in einem Klassenzimmer

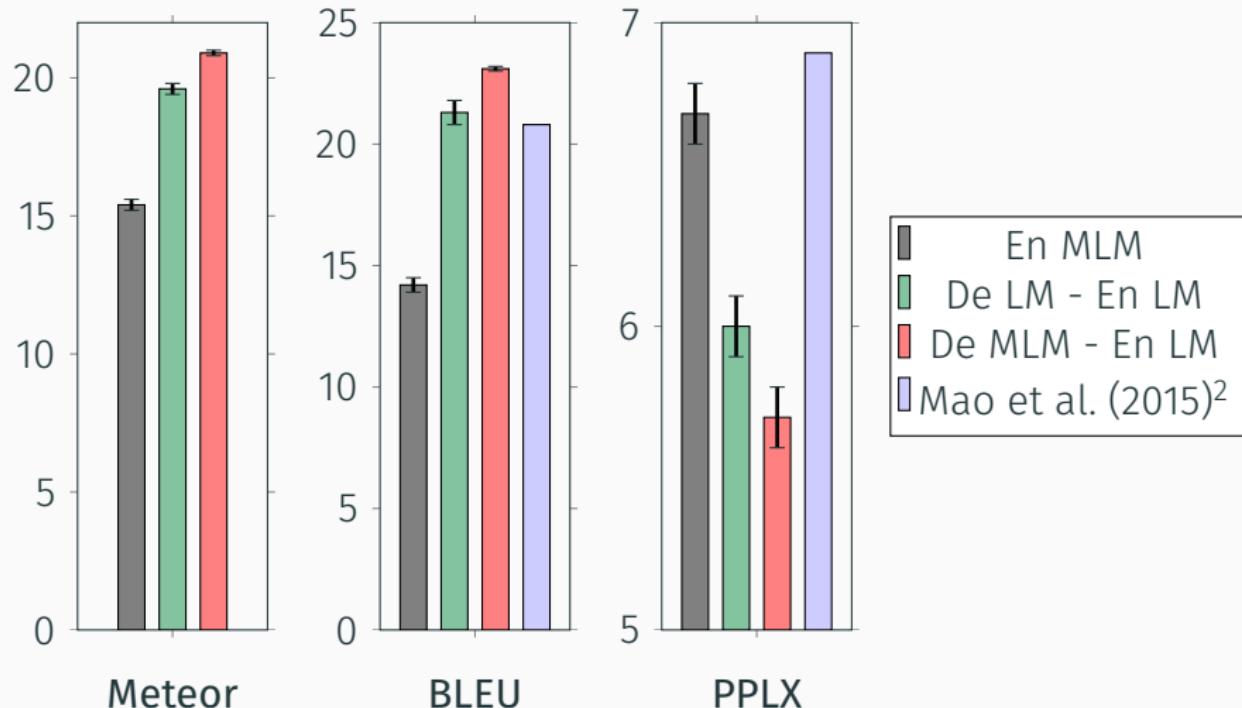


TARGET MODEL

## RESULTS

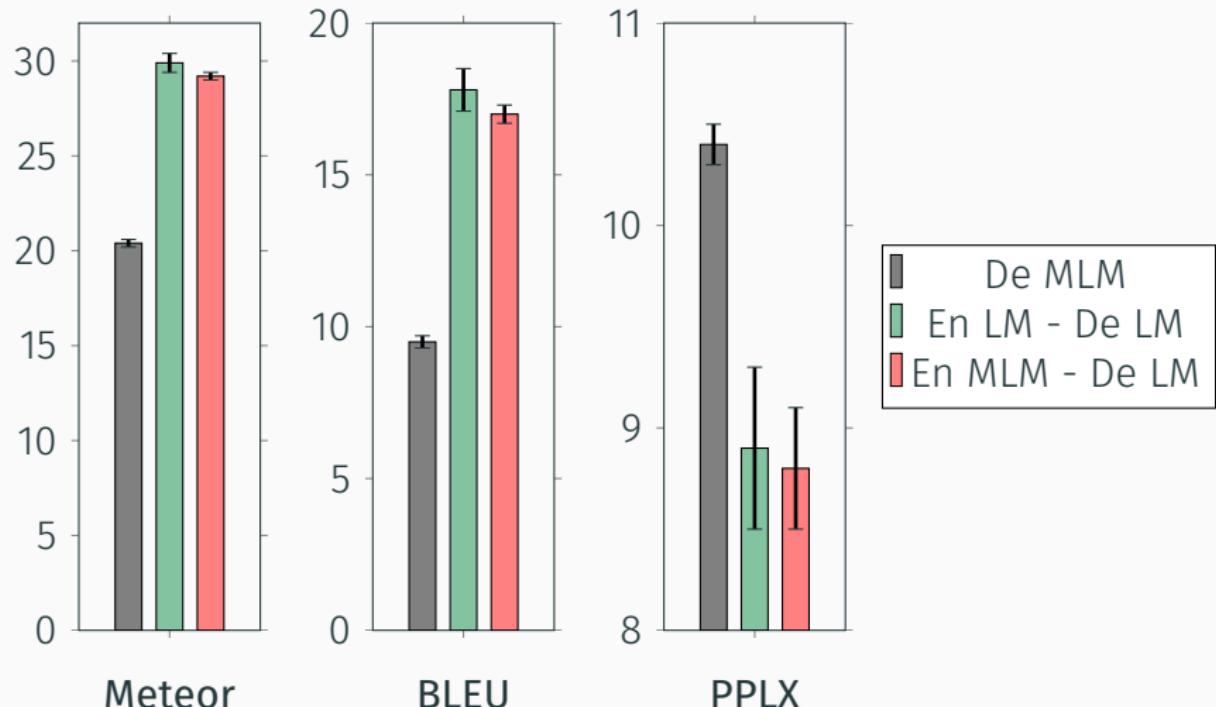
---

## MULTIMODAL TRANSLATION: ENGLISH RESULTS



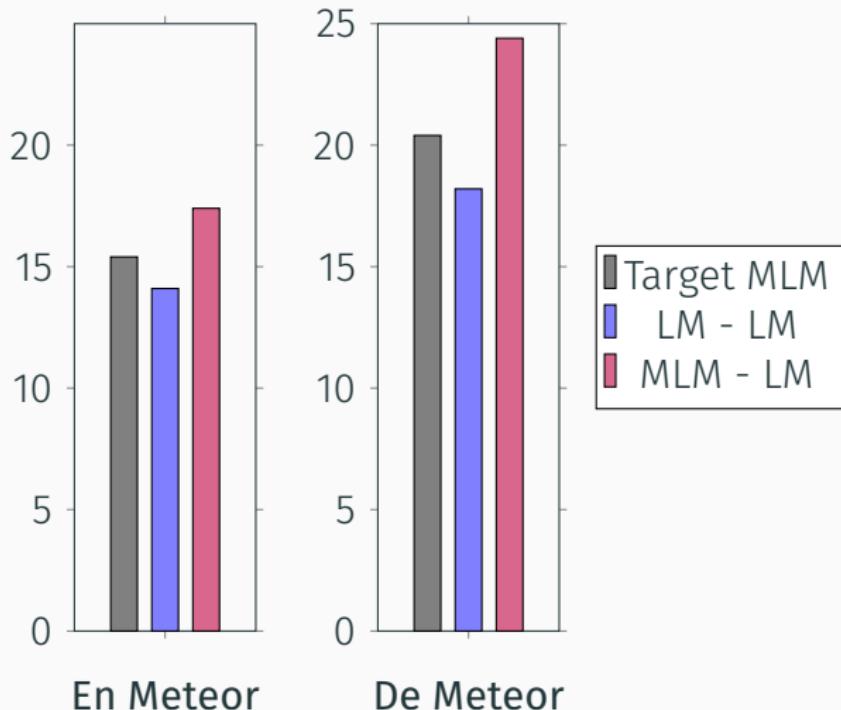
<sup>2</sup>Trained and evaluated over all references.

# MULTIMODAL TRANSLATION: GERMAN RESULTS<sup>3</sup>



<sup>3</sup>First non-English image description results.

# CROSSLINGUAL IMAGE DESCRIPTION RESULTS



Source descriptions automatically generated by Source-MLM

## EXAMPLE OF MULTIMODAL TRANSLATION



**MLM:** a man is standing on a grey rock in the foreground

**De Ref:** bergsteiger klettern auf einen sehr steilen eishang

**MLM-MLM:** tourists are climbing up a snowy slope

---

<sup>4</sup>Thousands of examples from all models at

<http://staff.fnwi.uva.nl/d.elliott/GroundedTranslation/>

## EXAMPLE OF CROSSLINGUAL IMAGE DESCRIPTION

De MLM:

ein mann und eine frau stehen  
an einem sandstrand mit dem  
meer im hintergrund



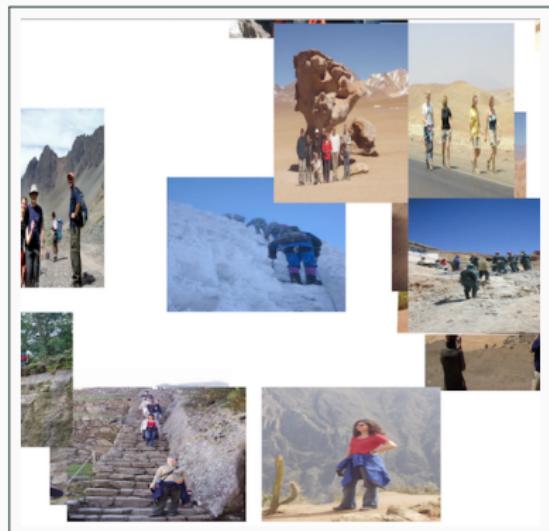
En MLM: a man with a black jacket and a black jacket is standing in a brown rocky desert landscape

En MLM-LM:

a man and a woman are  
standing in a reed boat on a  
lake

# VISUALISING THE EFFECT OF TRANSFERRING FEATURES

- t-SNE plots of the LSTM memory cell at  $t=0$
- MLM  $\rightarrow$  MLM: closer to pictures of snow!



(a) En MLM



(b) De MLM  $\rightarrow$  En MLM

- How well does this generalise to other languages?
- Attention-based Image Description [Xu et al., 2015]
- Compare with target-side translation retrieval with multimodal features [Hitschler and Riezler, 2016]
- Human judgements of generated descriptions
- Larger datasets (Shared Task at WMT16!)
- Multilingual video description, other tasks ...

## CONCLUSIONS

- Multilingual Image Description is a natural extension of Image Description
- MMLM transfers multimodal features between languages
- Transferring multilingual multimodal representations between languages improves image description quality
- Code: <http://github.com/elliottd/GroundedTranslation>

## APPENDICES

---

## COMPLETE ENGLISH RESULTS

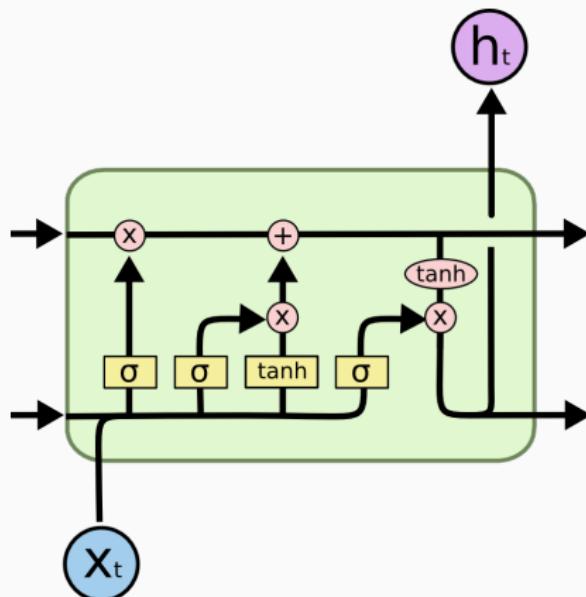
	BLEU4	Meteor	PPLX
En MLM	$14.2 \pm 0.3$	$15.4 \pm 0.2$	$6.7 \pm 0.0$
De LM → En LM	$21.3 \pm 0.5$	$19.6 \pm 0.2$	$6.0 \pm 0.1$
Mao et al. [2015]	20.8	—	6.92
De MLM → En MLM	$18.0 \pm 0.3$	$18.0 \pm 0.2$	$6.3 \pm 0.1$
De LM → En MLM	$17.3 \pm 0.5$	$17.6 \pm 0.5$	$6.3 \pm 0.0$
De MLM → En LM	<b><math>23.1 \pm 0.1</math></b>	<b><math>20.9 \pm 0.0</math></b>	$5.7 \pm 0.1$

## COMPLETE GERMAN RESULTS

	BLEU4	Meteor	PPLX
De MLM	$9.5 \pm 0.2$	$20.4 \pm 0.2$	$10.35 \pm 0.1$
En LM → De LM	<b><math>17.8 \pm 0.7</math></b>	<b><math>29.9 \pm 0.5</math></b>	$8.95 \pm 0.4$
En MLM → De MLM	$11.4 \pm 0.7$	$23.2 \pm 0.9$	$9.69 \pm 0.1$
En LM → De MLM	$12.1 \pm 0.5$	$24.0 \pm 0.3$	$10.2 \pm 0.7$
En MLM → De LM	<b><math>17.0 \pm 0.3</math></b>	<b><math>29.2 \pm 0.2</math></b>	$8.84 \pm 0.3$

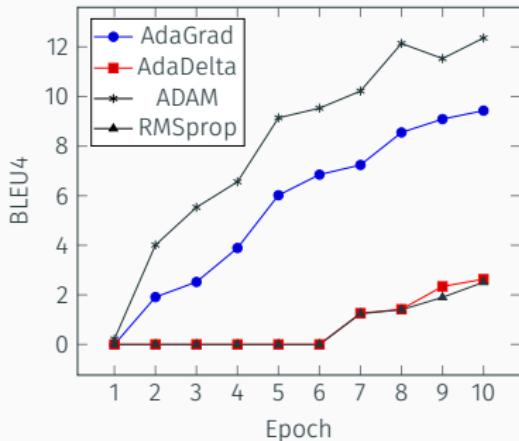
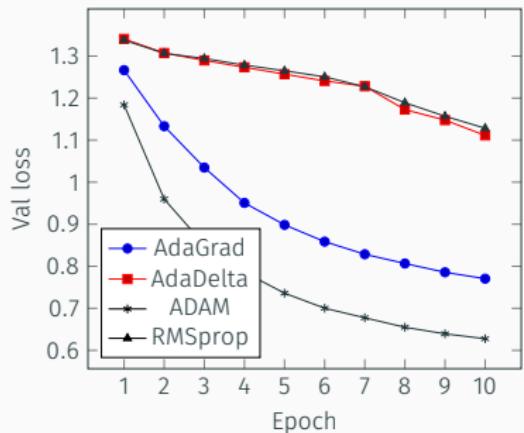
# RNN ARCHITECTURE: LONG-SHORT TERM MEMORY

[HOCHREITER AND SCHMIDHUBER, 1997]

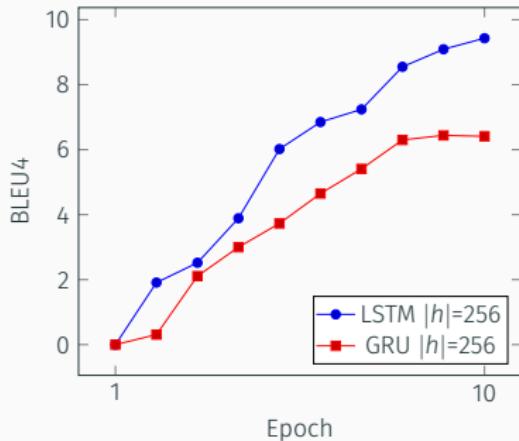
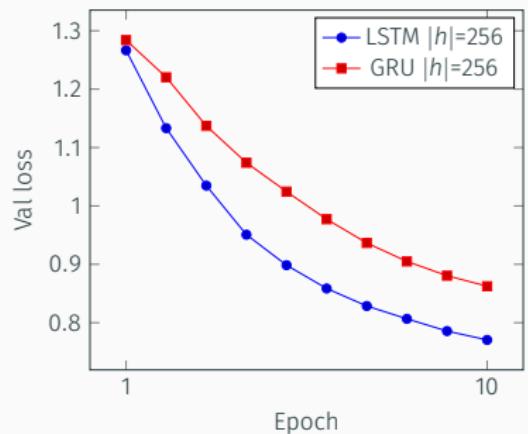


Credit: Christopher Olah

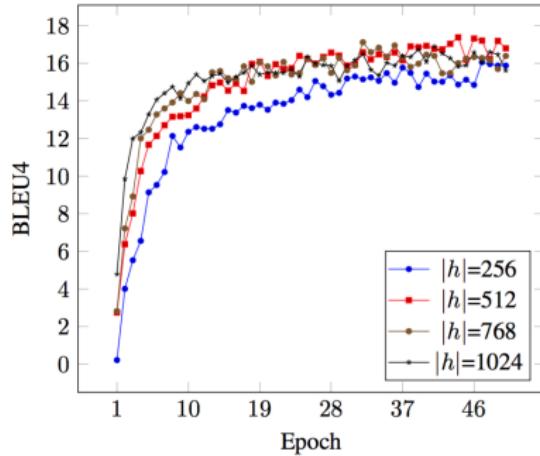
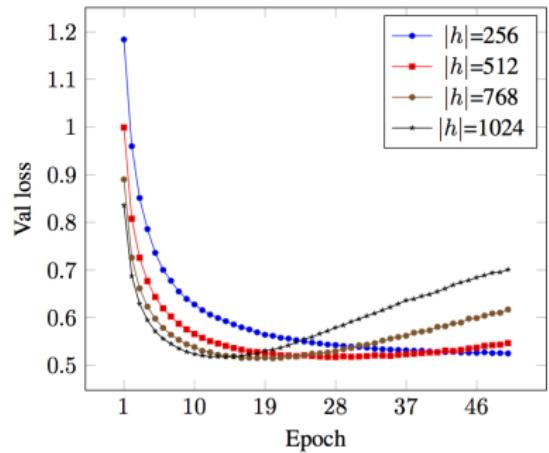
# EFFECT OF OPTIMISATION METHOD



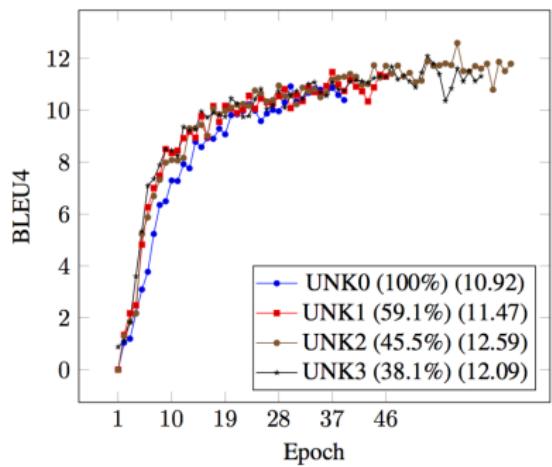
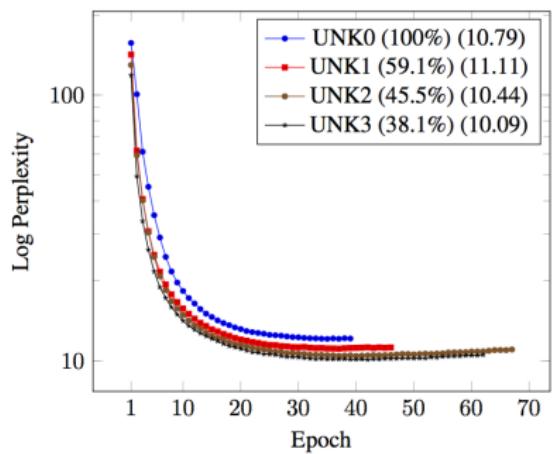
## EFFECT OF RNN TYPE



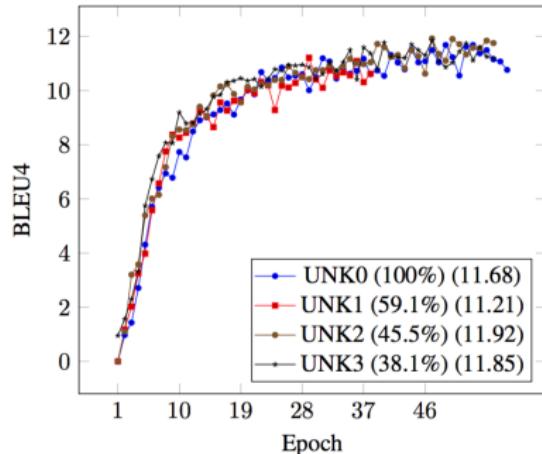
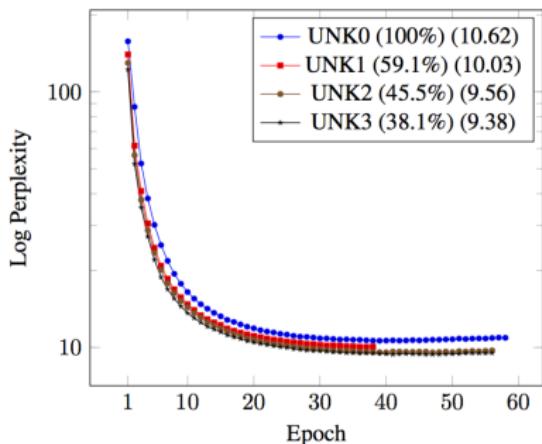
## EFFECT OF HIDDEN STATE SIZE



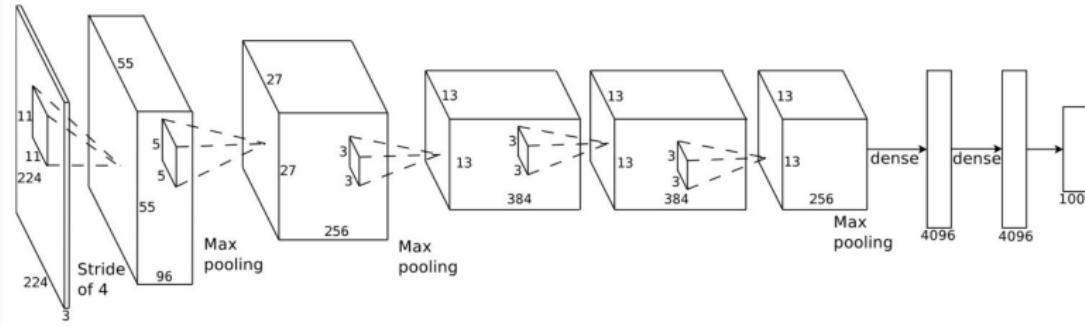
# EFFECT OF DECOMPOUNDING GERMAN WORDS



## EFFECT OF UNK THRESHOLD



# EXTRACTING CNN VISUAL FEATURES



Credit: Alex Krizhevsky

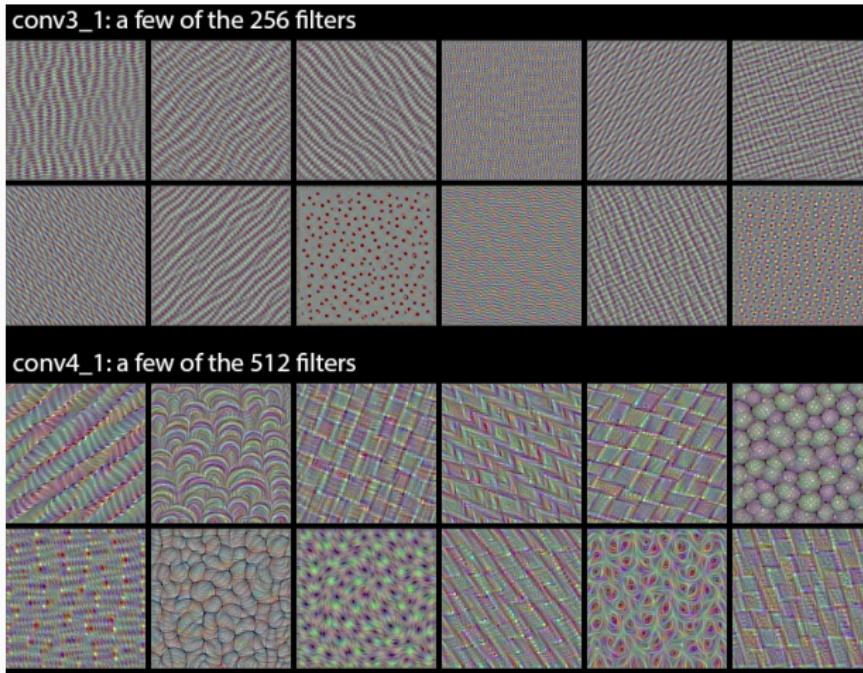
- Trained to predict 1000 object labels
- Over 1m training images
- Visual features transferred from the penultimate layer

# VISUALISING CNN FILTERS



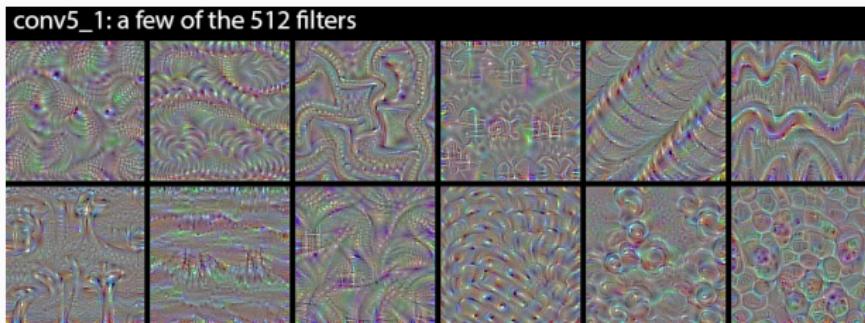
Credit: François Chollet

# VISUALISING CNN FILTERS



Credit: François Chollet

# VISUALISING CNN FILTERS



Credit: François Chollet

- Raffaella Bernardi, Ruken Cakici, Desmond Elliott, Aykut Erdem, Erkut Erdem, Nazli Ikizler-Cinbis, Frank Keller, Adrian Muscat, and Barbara Plank. Automatic description generation from images: A survey of models, datasets, and evaluation measures. *To appear in JAIR*, 2016.
- D Elliott, S Frank, and E Hasler. Multilingual image description with neural sequence models. *CoRR*, abs/1510.04709, 2015.
- Desmond Elliott and Frank Keller. Comparing Automatic Evaluation Measures for Image Description. In *ACL*, 2014.
- Ali Farhadi, M Hejrati, Mohammad Amin Sadeghi, P Young, C Rashtchian, J Hockenmaier, and David Forsyth. Every picture tells a story: Generating sentences from images. In *ECCV*, 2010.

## REFERENCES II

---

- Haoyuan Gao, Junhua Mao, Jie Zhou, Zhiheng Huang, Lei Wang, and Wei Xu. Are you talking to a machine? dataset and methods for multilingual image question answering. In *NIPS*, 2015.
- M. Grubinger, P. D. Clough, H. Muller, and D. Thomas. The IAPR TC-12 benchmark: A new evaluation resource for visual information systems. In *LREC*, 2006.
- J. Hitschler and S. Riezler. Multimodal Pivots for Image Caption Translation. *CoRR*, January 2016.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997.  
ISSN 0899-7667.

## REFERENCES III

---

- Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. *arXiv*, 2015.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- Junhua Mao, Wei Xu, Yi Yang, Jiang Wang, Zhiheng Huang, and Alan Yuille. Deep captioning with multimodal recurrent neural networks (m-RNN). *ICLR*, 2015.
- Carina Silberer and Mirella Lapata. Learning grounded meaning representations with autoencoders. In *Proceedings of ACL*, 2014.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR '15*, 2015.

## REFERENCES IV

- Ilya Sutskever, Oriol Vinyals, and Quoc V. V Le. Sequence to sequence learning with neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N.D. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 3104–3112. Curran Associates, Inc., 2014.
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, 2015.
- Subhashini Venugopalan, Marcus Rohrbach, Jeff Donahue, Raymond J. Mooney, Trevor Darrell, and Kate Saenko. Sequence to sequence - video to text. In *ICCV*, 2015.

## REFERENCES V

- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C. Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. *CoRR*, abs/1502.03044, 2015.