

How a General-Purpose Commonsense Ontology can Improve Performance of Learning-Based Image Retrieval

Rodrigo Toro Icarte[†] Jorge Baier^{‡,§} Cristian Ruz[‡] Alvaro Soto[‡]

[†]University of Toronto

[‡]Pontificia Universidad Católica de Chile

[§]Chilean Center for Semantic Web Research

IJCAI 2017

Motivation



Computer Vision

Simple ! Complex



Image classification: Kitchen

Computer Vision

Simple ! Complex



Image classification: Kitchen

Image captioning: “a woman in a chef coat holding bread loaves”

Computer Vision

Simple ! Complex



Image classification: Kitchen

Image captioning: “a woman in a chef coat holding bread loaves”

Image Q&A:

Q: What is the chef holding?

A: bread loaves

Computer Vision

Simple ! Complex



Image classification: Kitchen

Image captioning: “a woman in a chef coat holding bread loaves”

Image Q&A:

Q: What is the chef holding?

A: bread loaves

Computer Vision

Simple ! Complex

Learning everything from examples



Image classification: Kitchen

Image captioning: “a woman in a chef coat holding bread loaves”

Image Q&A:

Q: What is the chef holding?

A: bread loaves

Computer Vision

Simple ! Complex

Learning everything from examples

It does not scale (96.4% classification ! 32.2% captioning)



Image classification: Kitchen

Image captioning: “a woman in a chef coat holding bread loaves”

Image Q&A:

Q: What is the chef holding?

A: bread loaves

Computer Vision

Simple ! Complex

Learning everything from examples

It does not scale (96.4% classification ! 32.2% captioning)

Idea: Prior knowledge can fill the holes in our datasets.

Main research trends

Small hand-crafted ontologies

Free form text (e.g. Wikipedia)

Lexical ontologies (e.g. WordNet)

Main research trends

- Small hand-crafted ontologies

- Free form text (e.g. Wikipedia)

- Lexical ontologies (e.g. WordNet)

What about commonsense ontologies, such as ConceptNet?

ConceptNet (CN)

CN is a commonsense ontology.

ConceptNet (CN)

CN is a commonsense ontology.

Format

Concept₁ Relation type ! Concept₂

Relation types

AtLocation, HasProperty, IsA, SimilarSize, UsedFor, CapableOf, ...

ConceptNet (CN)

CN is a commonsense ontology.

Format

Concept₁ Relation type / Concept₂

Relation types

AtLocation, HasProperty, IsA, SimilarSize, UsedFor, CapableOf, ...

Examples

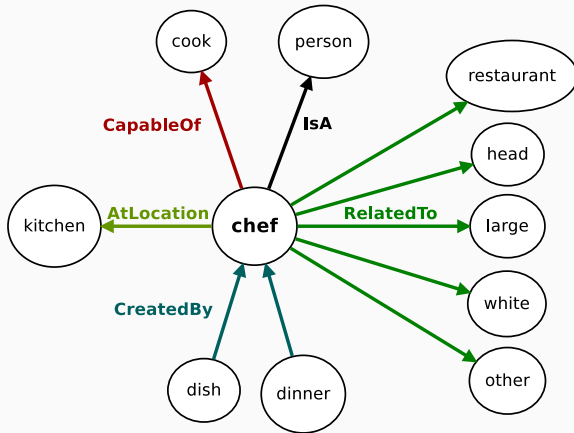
desk RelatedTo / office

computer AtLocation / office

office UsedFor / work

... and 8 million more

ConceptNet



CN is a great source of prior knowledge for Computer Vision.

CN is a great source of prior knowledge for Computer Vision.

- CN has millions of assertions (vs hand-crafted ontologies)

CN is a great source of prior knowledge for Computer Vision.

- CN has millions of assertions (vs hand-crafted ontologies)
- CN provides key knowledge to computers (vs Wikipedia)

CN is a great source of prior knowledge for Computer Vision.

- CN has millions of assertions (vs hand-crafted ontologies)
- CN provides key knowledge to computers (vs Wikipedia)
- CN is a rich source of commonsense knowledge (vs WordNet)

CN is a great source of prior knowledge for Computer Vision.

- CN has millions of assertions (vs hand-crafted ontologies)
- CN provides key knowledge to computers (vs Wikipedia)
- CN is a rich source of commonsense knowledge (vs WordNet)
- CN is simple to use (vs CYC)

Motivation

Motivation

Prior knowledge has a key role in Computer Vision.

Motivation

Prior knowledge has a key role in Computer Vision.

ConceptNet (CN) is a rich source of prior knowledge.

Motivation

Prior knowledge has a key role in Computer Vision.

ConceptNet (CN) is a rich source of prior knowledge.

Previous works

ConceptNet in Computer Vision

Task	w/o CN	w/ CN	CN gain
Image Tagging <small>Xie and He (2013)</small>	7.3%	7.6%	0.3%
Video Retrieval <small>de Boer et al. (2016)</small>	3.9%	3.1%	-0.8%
Image Riddles <small>Aditya et al. (2016)</small>	68.0%	68.7%	0.7%

More examples: Bicocchi et al. (2012), Le et al. (2013), others

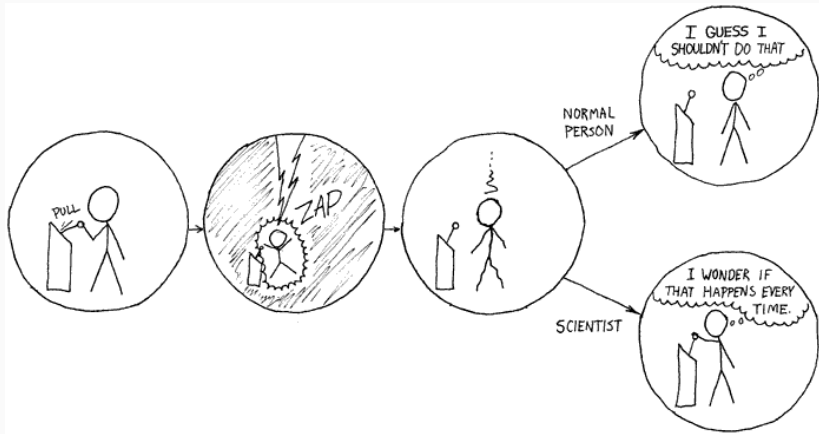
ConceptNet in Computer Vision

Task	w/o CN	w/ CN	CN gain
Image Tagging <small>Xie and He (2013)</small>	7.3%	7.6%	0.3%
Video Retrieval <small>de Boer et al. (2016)</small>	3.9%	3.1%	-0.8%
Image Riddles <small>Aditya et al. (2016)</small>	68.0%	68.7%	0.7%

More examples: Bicocchi et al. (2012), Le et al. (2013), others

... but we wanted to give CN another try.

... because we are scientists



Source: <https://xkcd.com/242/>

Motivation

Prior knowledge has a key role in Computer Vision.

ConceptNet (CN) is a rich source of prior knowledge.

Previous works

Motivation

Prior knowledge has a key role in Computer Vision.

ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.

Motivation

Prior knowledge has a key role in Computer Vision.

ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.

We don't care, we think CN is cool 🧐

Summary

Motivation

Prior knowledge has a key role in Computer Vision.

ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.

We don't care, we think CN is cool 🤓

Method

Summary

Motivation

Prior knowledge has a key role in Computer Vision.
ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.
We don't care, we think CN is cool 🧐

Method

CN for image retrieval...

Sentence Based Image Retrieval

“a woman in a chef coat holding bread loaves”



Sentence Based Image Retrieval

“a woman in a chef coat holding bread loaves”



Rank n images according to their *relevance* with respect to a natural language query.

Sentence Based Image Retrieval

“a woman in a chef coat holding bread loaves”



Rank n images according to their *relevance* with respect to a natural language query.



We used the 1000 Concept detectors trained by Fang et al.

We used the 1000 Concept detectors trained by Fang et al.



Prob	Concept
0.996	kitchen
0.920	preparing
0.800	food
0.796	cooking
0.590	making
...	...
0.236	woman

We used the 1000 Concept detectors trained by Fang et al.



Prob	Concept
0.996	kitchen
0.920	preparing
0.800	food
0.796	cooking
0.590	making
...	...
0.236	woman

t = “a woman in a chef coat holding bread loaves”

We used the 1000 Concept detectors trained by Fang et al.



Prob	Concept
0.996	kitchen
0.920	preparing
0.800	food
0.796	cooking
0.590	making
...	...
0.236	woman

$t =$ “a **woman** in a chef **coat holding bread** loaves”

$MIL(t, I) = P(\text{woman} | I) P(\text{coat} | I) P(\text{holding} | I) P(\text{bread} | I)$

t = \a woman in a chef coat holding bread loaves"

1 2 3 4 ... 512

MIL(t,l) = P(woman|l) P(coat|l) P(holding|l) P(bread|l)

$t = \backslash$ a woman in a chef coat holding bread loaves"

1 2 3 4 ... 512

$$\text{MIL}(t, I) = P(\text{woman}|I) \quad P(\text{coat}|I) \quad P(\text{holding}|I) \quad P(\text{bread}|I)$$

What are the limitations of this approach?

t = \a woman in a chef coat holding bread loaves"

1 2 3 4 ... 512

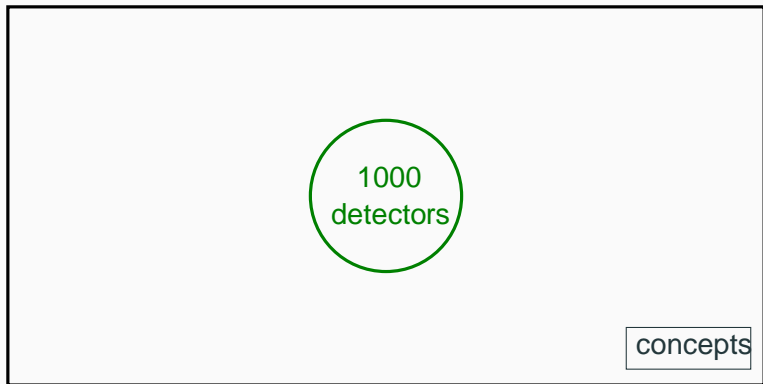
MIL(t, I) = P(woman|I) P(coat|I) P(holding|I) P(bread|I)

What are the limitations of this approach?

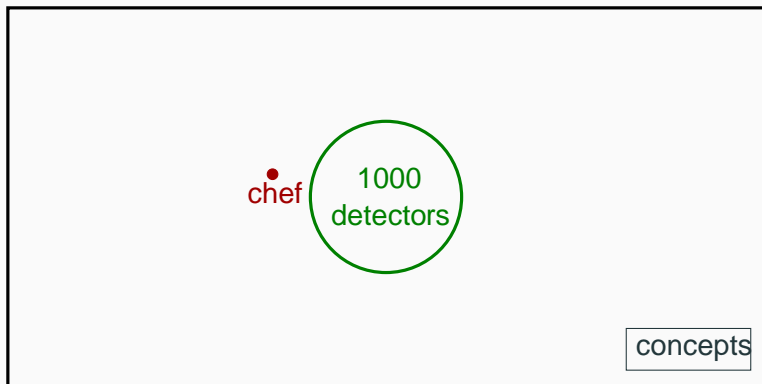
How can we detect a chef without a chef detector?



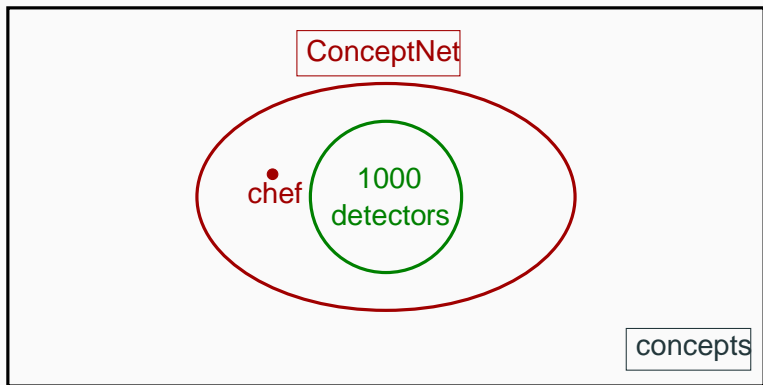
Idea: Augment the set of detectors using CN



Idea: Augment the set of detectors using CN



Idea: Augment the set of detectors using CN



Idea: Augment the set of detectors using CN

Word	Prob	Word	Prob
kitchen	0.996	dish	0.126
cook	0.796	white	0.091
restaurant	0.374	other	0.043
person	0.340	dinner	0.023
large	0.152	head	0.003

Word	Prob	Word	Prob
kitchen	0.996	dish	0.126
cook	0.796	white	0.091
restaurant	0.374	other	0.043
person	0.340	dinner	0.023
large	0.152	head	0.003

$$CN_{\text{MIN}}(\text{chef}) = 0.003$$

$$CN_{\text{AVG}}(\text{chef}) = 0.294$$

$$CN_{\text{MAX}}(\text{chef}) = 0.996$$

Database	r@1	r@5	r@10	median rank	mean rank
COCO 5K					
Baseline					
MIL	13.2	33.4	45.2	13	82.2
CN					
CN _{MIN}	12.2	31.4	43.4	15	77.0
CN _{AVG}	13.2	33.7	46.0	13	66.3
CN _{MAX}	12.2	32.1	44.1	14	73.0
CN Gain	0.0%	0.3%	0.8%	0	15.9

Summary

Motivation

Prior knowledge has a key role in Computer Vision.

ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.

We don't care, we think CN is cool

Method

CN for image retrieval...

Summary

Motivation

Prior knowledge has a key role in Computer Vision.

ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.

We don't care, we think CN is cool

Method

CN for image retrieval...sucks!

$P_c(\text{chef}|\text{l})$

$$P_c(\text{chef} | I) = P(\text{chef} | \text{cook}, I)P(\text{cook} | I) + P(\text{chef} | \text{: cook}, I)P(\text{: cook} | I)$$

$$P_c(\text{chef} | I) = P(\text{chef} | \text{cook}, I)P(\text{cook} | I) + P(\text{chef} | : \text{cook}, I)P(: \text{cook} | I)$$

$$P_c(\text{chef} | I) = P(\text{chef} | \text{cook}, I) \cdot 0.796 + P(\text{chef} | : \text{cook}, I) \cdot 0.204$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}|\text{l})P(\text{cook}|\text{l}) + P(\text{chef}|\text{: cook}|\text{l})P(\text{: cook}|\text{l})$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}|\text{l}) \cdot 0.796 + P(\text{chef}|\text{: cook}|\text{l}) \cdot 0.204$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}) \cdot 0.796 + P(\text{chef}|\text{: cook}) \cdot 0.204$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}|\text{l})P(\text{cook}|\text{l}) + P(\text{chef}|\text{: cook}|\text{l})P(\text{: cook}|\text{l})$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}|\text{l}) \ 0:796 + P(\text{chef}|\text{: cook}|\text{l}) \ 0:204$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}) \ 0:796 + P(\text{chef}|\text{: cook}) \ 0:204$$

$$P_c(\text{chef}|\text{l}) = 0:1413 \ 0:796 + 0:0003 \ 0:204$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}|\text{l})P(\text{cook}|\text{l}) + P(\text{chef}|\text{: cook}|\text{l})P(\text{: cook}|\text{l})$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}|\text{l}) \ 0:796 + P(\text{chef}|\text{: cook}|\text{l}) \ 0:204$$

$$P_c(\text{chef}|\text{l}) = P(\text{chef}|\text{cook}) \ 0:796 + P(\text{chef}|\text{: cook}) \ 0:204$$

$$P_c(\text{chef}|\text{l}) = 0:1413 \ 0:796 + 0:0003 \ 0:204$$

$$P_c(\text{chef}|\text{l}) = 0:112$$

$P_{\text{cook}}(\text{chef} I)$	0:1125
$P_{\text{kitchen}}(\text{chef} I)$	0:0549
$P_{\text{dish}}(\text{chef} I)$	0:0016
$P_{\text{person}}(\text{chef} I)$	0:0011
$P_{\text{dinner}}(\text{chef} I)$	0:0011
$P_{\text{head}}(\text{chef} I)$	0:0009
$P_{\text{other}}(\text{chef} I)$	0:0009
$P_{\text{white}}(\text{chef} I)$	0:0006

$P_{\text{cook}}(\text{chef} I)$	0:1125
$P_{\text{kitchen}}(\text{chef} I)$	0:0549
$P_{\text{dish}}(\text{chef} I)$	0:0016
$P_{\text{person}}(\text{chef} I)$	0:0011
$P_{\text{dinner}}(\text{chef} I)$	0:0011
$P_{\text{head}}(\text{chef} I)$	0:0009
$P_{\text{other}}(\text{chef} I)$	0:0009
$P_{\text{white}}(\text{chef} I)$	0:0006

$CNE_{\text{MIN}}(\text{chef})$	= 0.0006
$CNE_{\text{AVG}}(\text{chef})$	= 0.0217
$CNE_{\text{MAX}}(\text{chef})$	= 0.1125

CN + ESPGAME Score

Database	r@1	r@5	r@10	median rank	mean rank
COCO 5K					
Baseline					
MIL	13.2	33.4	45.2	13	82.2
CN + ESPGAME					
CN _{E_{MIN}}	14.3	34.6	46.6	12	68.3
CN _{E_{AVG}}	14.6	35.6	48.0	12	61.2
CN _{E_{MAX}}	14.3	35.9	48.2	12	60.6
CN Gain	1.4%	2.5%	3.0%	1	21.6

t = \a woman in a chef coat holding bread loaves"

1

2

3

4

...

512

$t =$ “a woman in a chef coat holding bread loaves”

1 2 3 4 ... **512**

1 2 3 4 ... **35**

$t =$ “those bagels are plain with nothing on them”

1 2 3 4 ... **360**

$t =$ “those bagels are plain with nothing on them”

1	2	3	4	...	360
---	---	---	---	-----	------------

1	2	3	4	...	2
---	---	---	---	-----	----------

Summary

Motivation

Prior knowledge has a key role in Computer Vision.
ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.
We don't care, we think CN is cool 🧐

Method

CN for image retrieval... sucks!

Motivation

Prior knowledge has a key role in Computer Vision.
ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.
We don't care, we think CN is cool 🤓

Method

CN for image retrieval... sucks!
CN + ESPGAME for image retrieval... works!

Summary

Motivation

Prior knowledge has a key role in Computer Vision.
ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.
We don't care, we think CN is cool 🧐

Method

CN for image retrieval... sucks!
CN + ESPGAME for image retrieval... works!

Contribution

Results and Discussion

Database	r@1	r@5	r@10	median rank	mean rank
COCO 5K					
NeuralTalk (Vinyals et al., 2015)	6.9	22.1	33.6	22	72.2
GMM+HGLMM (Klein et al., 2015)	10.8	28.3	40.1	17	49.3
BRNN (Karpathy and Fei-Fei, 2015)	10.7	29.6	42.2	14	–
MIL (our baseline)	15.7	37.8	50.5	10	53.6
CNE_{MAX} (our method)	16.2	39.1	51.9	10	44.4
LVQ (Lin and Parikh, 2016)	16.7	40.5	53.8	–	–
OE (Vendrov et al., 2016)	18.0	–	57.6	7.0	35.9

higher is better

lower is better

Results and Discussion

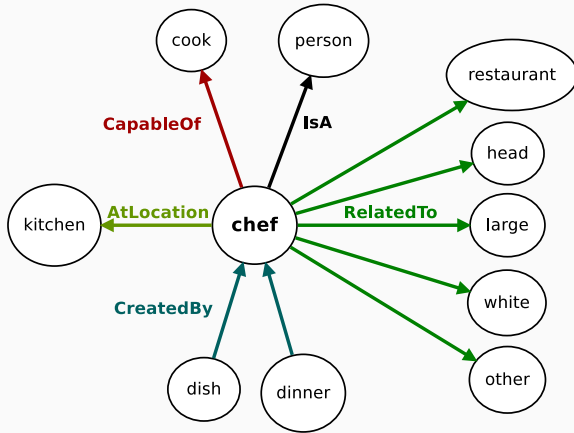
Database	r@1	r@5	r@10	median	mean
COCO 5K				rank	rank
MIL (our baseline)	13.2	33.4	45.2	13	82.2
CNE _{MAX} (our method)	14.3	35.9	48.2	12	60.6
CN Gain	1.1%	2.5%	3.0%	1	21.6

Results and Discussion

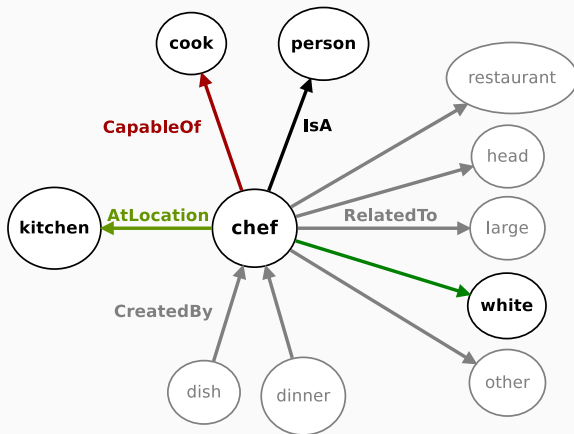
Database	r@1	r@5	r@10	median	mean
COCO 5K				rank	rank
MIL (our baseline)	13.2	33.4	45.2	13	82.2
CNE _{MAX} (our method)	14.3	35.9	48.2	12	60.6
CN Gain	1.1%	2.5%	3.0%	1	21.6

Why is CN helping this time?

Why is CN helping this time?



Why is CN helping this time?



Summary

Motivation

Prior knowledge has a key role in Computer Vision.
ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.
We don't care, we think CN is cool 🤓

Method

CN for image retrieval... sucks!
CN + ESPGAME for image retrieval... works!

Contribution

Motivation

Prior knowledge has a key role in Computer Vision.
ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.
We don't care, we think CN is cool 🤓

Method

CN for image retrieval... sucks!
CN + ESPGAME for image retrieval... works!

Contribution

We can exploit commonsense ontologies in Computer Vision

Motivation

Prior knowledge has a key role in Computer Vision.
ConceptNet (CN) is a rich source of prior knowledge.

Previous works

They suggest that CN sucks.
We don't care, we think CN is cool 🤓

Method

CN for image retrieval... sucks!
CN + ESPGAME for image retrieval... works!

Contribution

We can exploit commonsense ontologies in Computer Vision,
but this knowledge must be filtered in a meaningful way (e.g.
using ESPGAME).

Acknowledgements



Becas Chile — Magister en el Extranjero



FONDECYT
Fondo Nacional de Desarrollo
Científico y Tecnológico

FONDECYT 1151018 and 1150328



UNIVERSITY OF
TORONTO

School of
Graduate Studies

UofT School of Graduate Studies



CANADIAN ARTIFICIAL INTELLIGENCE ASSOCIATION
ASSOCIATION POUR L'INTELLIGENCE ARTIFICIELLE AU CANADA

Canadian Artificial Intelligence Association



Our code: <https://bitbucket.org/RTorolcarte/cn-detectors>

Our code: <https://bitbucket.org/RTorolcarte/cn-detectors>

If you want to share ideas about commonsense knowledge in Computer Vision, please come to check my poster :)

Our code: <https://bitbucket.org/RTorolcarte/cn-detectors>

If you want to share ideas about commonsense knowledge in Computer Vision, please come to check my poster :)

Thank you!

References I

- Somak Aditya, Yezhou Yang, Chitta Baral, and Yiannis Aloimonos. Answering image riddles using vision and reasoning through probabilistic soft logic. *arXiv preprint arXiv:1611.05896*, 2016.
- Nicola Bicocchi, Matteo Lasagni, and Franco Zambonelli. Bridging vision and commonsense for multimodal situation recognition in pervasive systems. In *Pervasive Computing and Communications (PerCom), 2012 IEEE International Conference on*, pages 48–56. IEEE, 2012.
- Maaïke de Boer, Klamer Schutte, and Wessel Kraaij. Knowledge based query expansion in complex multimedia event detection. *Multimedia Tools and Applications*, 75(15):9025–9043, 2016.
- Hao Fang, Saurabh Gupta, Forrest Iandola, Rupesh K Srivastava, Li Deng, Piotr Dollár, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C Platt, et al. From captions to visual concepts and back. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1473–1482, 2015.
- Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.
- Benjamin Klein, Guy Lev, Gil Sadeh, and Lior Wolf. Associating neural word embeddings with deep image representations using fisher vectors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4437–4446, 2015.
- Dieu-Thu Le, Jasper RR Uijlings, and Raffaella Bernardi. Exploiting language models for visual recognition. In *EMNLP*, pages 769–779, 2013.
- Xiao Lin and Devi Parikh. Leveraging visual question answering for image-caption ranking. In *European Conference on Computer Vision*, pages 261–277. Springer, 2016.
- Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. Order-embeddings of images and language. In *ICLR*, 2016.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3156–3164, 2015.
- Lexing Xie and Xuming He. Picture tags and world knowledge: learning tag relations from visual semantic sources. In *Proceedings of the 21st ACM international conference on Multimedia*, pages 967–976. ACM, 2013.