

Using linguistic prior knowledge for deep learning

Dr Nina Dethlefs
n.dethlefs@hull.ac.uk

14 Feb 2018

Seminar: Bangor University



UNIVERSITY
OF HULL

Natural Language Generation (NLG)

Generating natural language from non-linguistic representations.

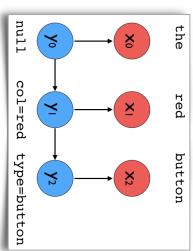
UNIVERSITY
OF HULL

Building NLG Systems

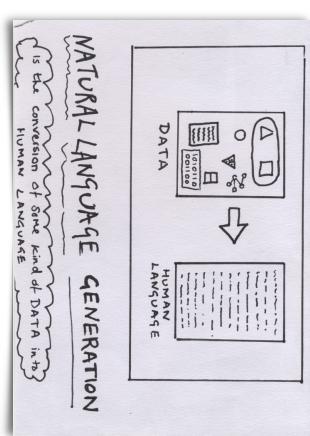
3 Approaches.

Mapping rules: human experts help system developers design sophisticated rules to generate text from data in a specific domain.

Statistical approaches: learn a mapping from aligned data or (more recently) from unaligned data.

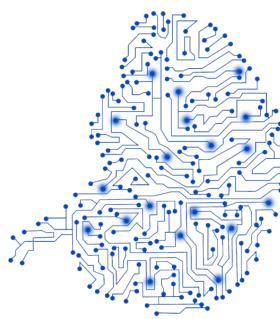


Deep learning: learn a direct mapping from inputs to outputs.



<https://www.flickr.com/photos/triptech/7690700352>

Deep Learning

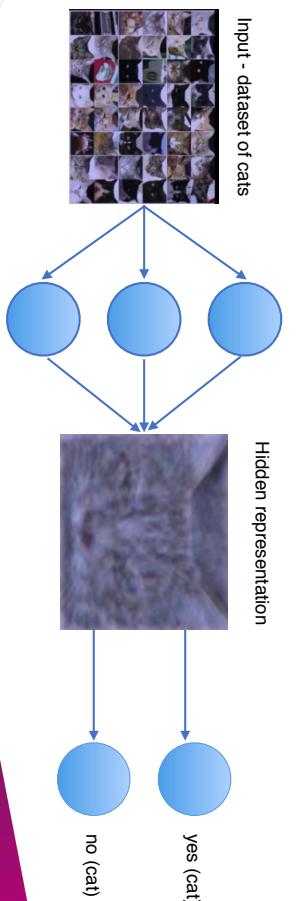


Models learn increasingly **abstract** representations of inputs - at different layers of abstraction - and recognise patterns in unstructured data.

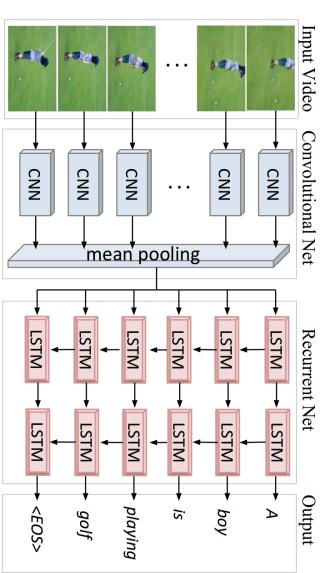
Applications: self-driving cars, natural language processing, image/video recognition, medical imaging, stock market prediction, ...

Computationally intensive.

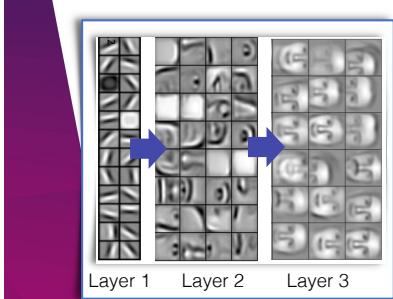
Imagine processing - abstract features



Video caption generation



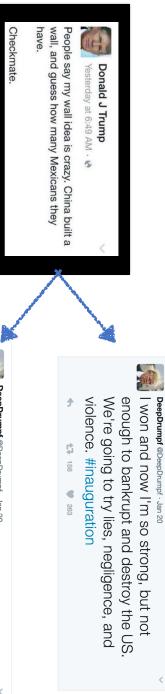
Deep Learning - Basic ideas



Sequence generation



Sequence generation



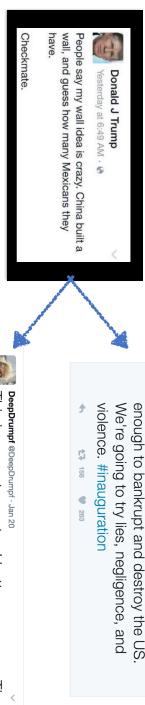
Microsoft's Tay



MIT's DeepDrumpf

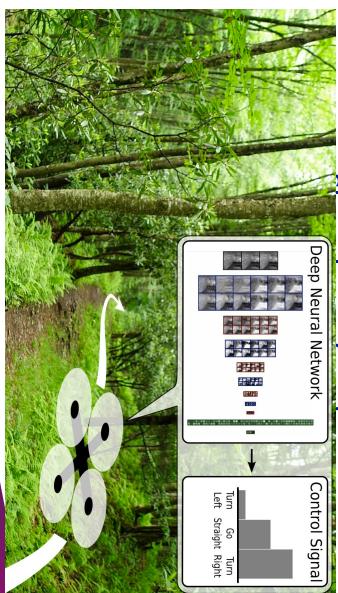
<https://twitter.com/deepdrumpf?lang=en>

Sequence generation



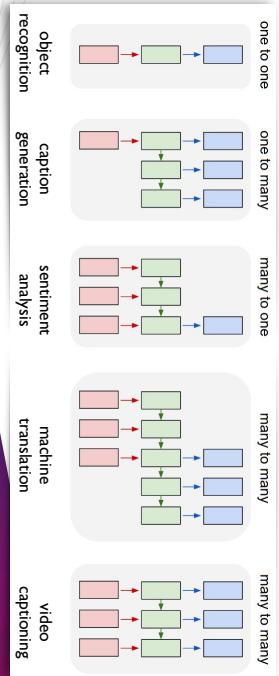
Flying drones

Drones learn to fly unsupervised in new environments.



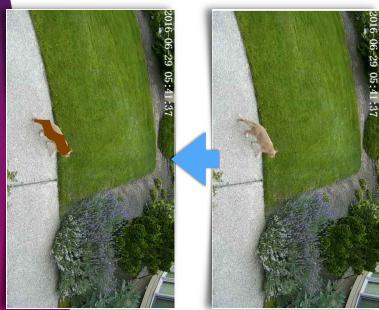
The “best” neural network

Different types of neural nets for different tasks...



Deep learning - summary

- Deep learning works - given the right setup and conditions - and can be applied to many useful tasks
- Different neural networks are suitable for different tasks, including language tasks
- Lots of parameters and space for improvement...



Domain Transfer for Natural Language Generation



“Folk wisdoms” of deep learning

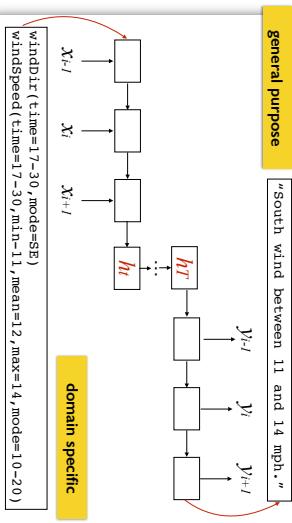
- To train good models, we need a lot of data
- Data need to be “balanced”
- Very general rules of thumb for parameter tuning
- Generally run on a GPU and / or in a distributed way
- (Don’t expect to understand what the model learnt and why)



Deep learning for NLG

Learn a conditional mapping from inputs to outputs.

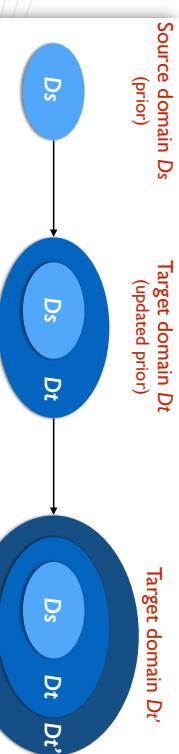
The Encoder-Decoder LSTM:
encodes and input sequence into
an abstract hidden representation
and learns to decode it to a
(different) output sequence.



Motivation

Develop **NLG systems** from datasets relatively quickly (also) for new / flexible / dynamic domains.

Idea of this work: instead of collecting ever more data to make deep learning work, can we structure datasets differently to make the most of them?



Domain Adaptation in NLP

Not a new idea in other areas of NLP: machine translation, NER,
capitalisation, shallow parsing.

Approaches often reuse prior knowledge.

Transfer probs from a source to a target domain, and then **adapt** them
(Chelba and Accero, 2006; Daume 2007, many more since)

Normally rely on a common feature representation.

Recent work on NLG across domains using **synthetic data** and shared input representations (Wen et al., 2016).

(Part of) The Problem with NLG

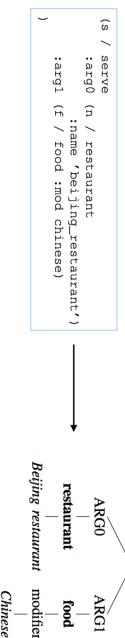
- 1a. Beijing restaurant is serving Chinese food.
- 1b. Beijing restaurant is a nice place. It serves Chinese food.
- 1c. Beijing restaurant is a Chinese restaurant.

*different ways of saying
the same thing?*

Flat semantic representation (for realizations 1a-1c)

inform(name='beijing_restaurant', food=chinese)

AMR and corresponding tree representation (for realization 1a.)



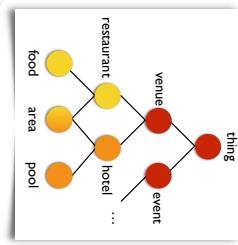
(Part of) The Problem with NLG

Underspecified semantic forms may not just pretend that these are the same:

1. Beijing restaurant serves Chinese food.
2. Beijing restaurant is a nice place. It serves Chinese food.
3. Beijing restaurant is a Chinese restaurant.

but also that they're not nothing to with:

4. The blue cube touches the red sphere.



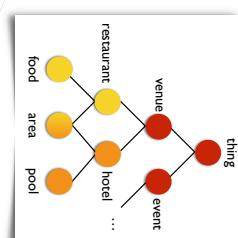
(Part of) The Problem with NLG

Underspecified semantic forms may not just pretend that these are the same:

1. Beijing restaurant serves Chinese food.
2. Beijing restaurant is a nice place. It serves Chinese food.
3. Beijing restaurant is a Chinese restaurant.

but also that they're not nothing to with:

4. The blue cube touches the red sphere



Learning Constructions

Constructions in cognitive grammar are [conventionalised form-meaning pairs](#) at different levels of abstraction: (idioms), partially lexically filled, fully general phrases.

1. Elmo the car gopping. (SOV)

2. Dacking Elmo the car. (VSO)

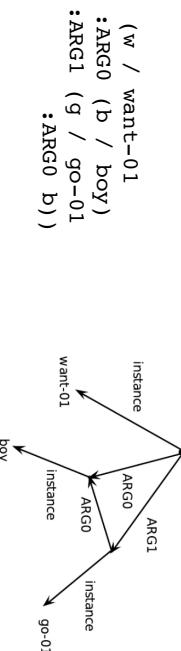
3. Elmo blicking the car. (SVO)

Can we use constructions in our **x** to **y** mappings?



Goldberg 2005

Abstract Meaning Representations



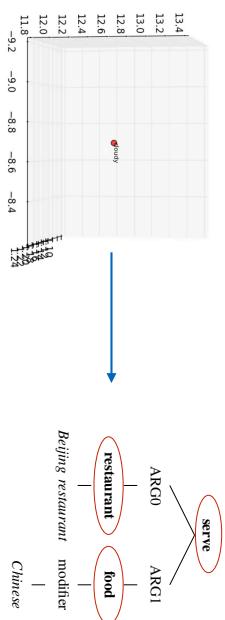
Semantic representations as directed graphs — representing a shared set of [semantic roles](#), [named entities](#), [modality](#), [negation](#), [questions](#), [quantities](#), etc.

Dealing with words

AMRs help us find structural similarities but don't solve the "unseen word problem"—for example, to a neural net, these are complete different:

The Kirin	is	a child-friendly restaurant
The Kirin	is	a dog-friendly restaurant
The Kirin	is	an exquisite restaurant

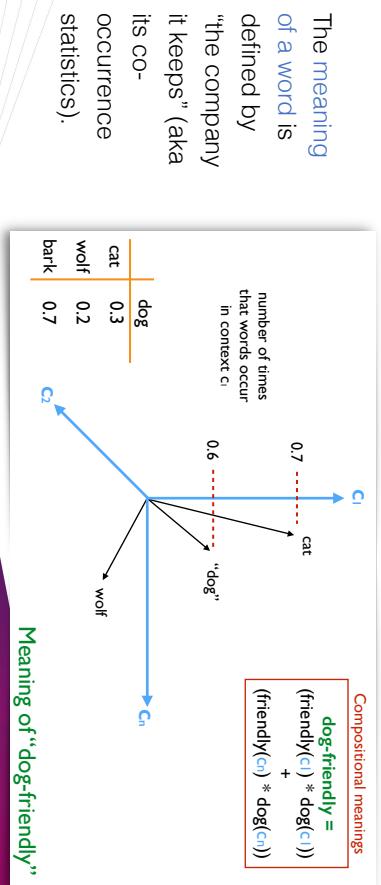
 UNIVERSITY OF HULL



Semantic classes have the same identifier if they occur in the same context.

© 2018 UNIVERSITY OF HULL

Hypothesis: using AMRs + deep learning, we can learn enough lexical-syntactic patterns to generate language in unseen target domains.



Source and target domains

Image recognition

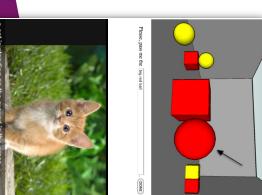
Image recognition
(caption generation)

Navigation
(planning)

Tourist info
(dialogue)



Source domain: GRE3D7 objects in a spatial arrangement. Referring expressions and spatial relations “The blue cube next to the yellow sphere”.



Navigation (planning)

Navigation
(planning)



Source domain: GIVE2 navigation in a spatial virtual environment “go left through the door and press the red button”. Referring expressions, spatial language, imperatives.

Target domain: SAIL navigation in real environments. Similar but more varied.

<https://www.youtube.com/watch?v=5-vsSnFX6kI>

Tourist info
(dialogue)



Source domain: SFXRestaurants spoken dialogue system utterances enquiring about user preferences and making recommendations.

Target domain: SFXHotels spoken dialogue system in a similar domain but using different slots and slightly different distribution and sample size.

<https://www.youtube.com/watch?v=EZw08xv2HEl&t=23s>

Datasets - statistics

Dataset	Examples	Vocab. size	Ave. len.	NPs	SRs	Trans. Cl.	Intr. Cl.	Rel. Cl.	Imperatives	Conj.
GRE	4,480	195	3.50	4,969	1,084	45	0	15	3	2,329
REFCoco	142,051	10,046	3.50	302,703	54,484	384	6,634	794	229	
GIVE	1,756	467	3.30	706	1,716	294	17	4	1,043	100
SAIL	816	295	7.95	598	895	218	64	5	656	135
SFR	6,198	2,058	12.91	6,094	956	2,597	872	598	3	905
SFH	6,384	1,061	12.13	6,403	1,192	2,247	938	761	1	830

Similar linguistic patterns in datasets but with different distributions.

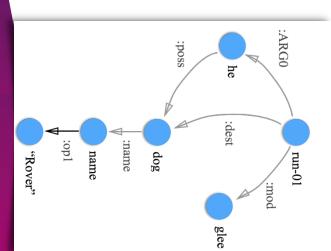
Experimental setup

Three experiments:

In-domain training: train and test in the same domain.

Out-of-domain training: train in one domain, test in another.

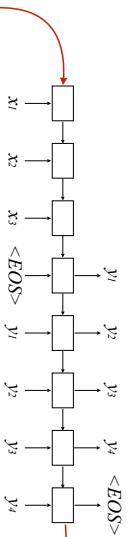
Cross-domain training: train in target domain, but use prior knowledge from a different domain.



(Deep) Learning from AMRs

Improve performance of NLG by leveraging general-purpose linguistic knowledge.

Beijing restaurant serves Chinese food.



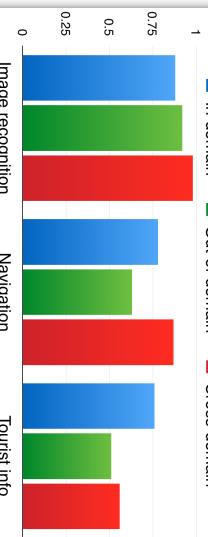
```
(s / serve
  :arg0 (n / restaurant
    :name 'beijing_restaurant')
  :arg1 (f / food :mod chinese)
)
```

Trained 4-layered LSTMs with 50 hidden units, batch size 12 over 10,000 training epochs.

Case study: Results

Image recognition

■ In-domain ■ Out-of-domain ■ Cross-domain



- Out-of-domain: 75%-100% of original performance
- Cross-domain: up to 110% of original performance.

Results: Linguistic Patterns

Complex NP with spatial relation and temporal adverb

```
(e1 / obj :time (n / now)
:domain (b1 / obj :mod property
:location (l / obj
:opl ( n / on [near, by])) )
```

"Now, the blue circle on the green sphere."

"Now, the green button by the window."

"Now, an Indian restaurant near Pacific Heights."

Results: Linguistic Patterns

Imperative clause construction (with relative clause)

```
(e1 / event :arg0 (y / you)
:arg1 (b1 / obj :mod property
:location (w / obj :opl (o1 / on [in]) )
:mode imperative )
```

"Click the red button (that is) on the wall."

"Try Chinese restaurant Kirin in the Pacific Heights area."

Results: Linguistic Patterns

Transitive clause construction

```
(e1 / event :arg0 (b1 / obj :mod property)
:arg1 (b2 / obj :mod property))
```

"The yellow sphere (that is) touching the blue box."

"Source restaurant serves Italian food."

Results: similarity and error

Objective evaluation using
BLEU-3 and -4 score (n-gram

[similarity](#) with human
examples) and [semantic](#)

error.

No significant differences.

	GRE-HUMAN	REFCOCO	GIVE	SAIL	SFXR
in-domain	0.90	0.88	0.82	0.02	0.016
GRE trained from GIVE	0.79	0.78	0.77	0.09	
REFCoco trained from GRE	0.82	0.81	0.76	0.10	0.04
GIVE trained from GRE	0.75	0.69	0.69	0.08	
SAIL trained from GIVE	0.94	0.92	0.04		
SFXR trained from GIVE	0.91	0.84	0.02		
GRE trained from SFXR	0.23	0.16	0.29		
REFCoco trained from SFXR	0.68	0.63	0.10		
GIVE trained from SFXR	0.61	0.51	0.12		
SAIL trained from SFXR	0.68	0.63	0.10		
GIVE with GIVE prior	0.99	0.98	0.0		
REFCoco with GRE prior	0.96	0.77	0.01		
GIVE with GRE prior	0.89	0.87	0.06		
SAIL with GIVE prior	0.80	0.72	0.05		
SFXR with SFXR prior	0.74	0.56	0.08		

Results: naturalness

Subjective evaluation with 204 human judges from AMT, rating 3425 utterances.

"The utterance is natural i.e. could have been produced by a human"

Prior knowledge is better than pure domain transfer. No other significant differences.

	SYSTEM	NATURALNESS
Human	GRE-HUMAN	2.97 (4)
	REFCOCO	3.43 (4)
	GIVE	3.77 (4)
	SAIL	4.36 (5)
	SFXR	4.17 (4)
SFXH		3.93 (4)
	GRE-HUMAN	3.12 (3)
	REFCOCO	3.45 (4)
	GIVE	2.76 (2)
	SAIL	4.02 (4)
	SFXR	3.67 (4)
SFXH		3.95 (4)
	GRE trained from GIVE	3.11 (3)
	REFCOCO trained from GRE	3.20 (3)
	GIVE trained from Give	2.03 (2)
	SAIL trained from Give	2.95 (3)
	SFXH trained from SFXR	3.45 (3)
	GRE with GIVE prior	3.29 (3)
	REFCOCO with GRE prior	3.41 (4) *
	GIVE with GRE prior	3.09 (3) *
prior	SAIL with GIVE prior	3.44 (4) *
	SFXH with SFXR prior	3.58 (4)

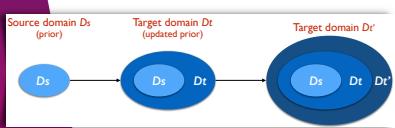
Did it work?

Yes:
Model learns lexical-syntactic patterns from AMRs that are not domain-specific

Out-of-domain results are **acceptable** in some domains, similar to in-domain training in others

Learning from source domain **priors** gives us up to **10% extra performance**

Limitations:
AMR annotation is an overhead.



Conclusions

Deep learning works well, but only with enough data of the right kind.

More "intelligence" needed in algorithms — transfer of prior knowledge across tasks.

More control needed to constrain decisions — linguistic knowledge is one way to achieve this for NLP tasks.

(Synthetic data seems a bad idea)

