



Learning to Compose Neural Networks for Question Answering (a.k.a. Dynamic Neural Module Networks)

Authors: Jacob Andreas, Marcus Rohrbach, Trevor Darrell, Dan Klein
Presented by: K.R. Zentner



Basic Outline

- Problem statement
- Brief review of Neural Module Networks
- New modules
- Learned layout predictor
- Some minor additions
- Results
- Conclusion



Problem Statement

Would like to have a single algorithm for a variety of question answering domains.

More precisely, given a question q and a world w , produce an answer y .

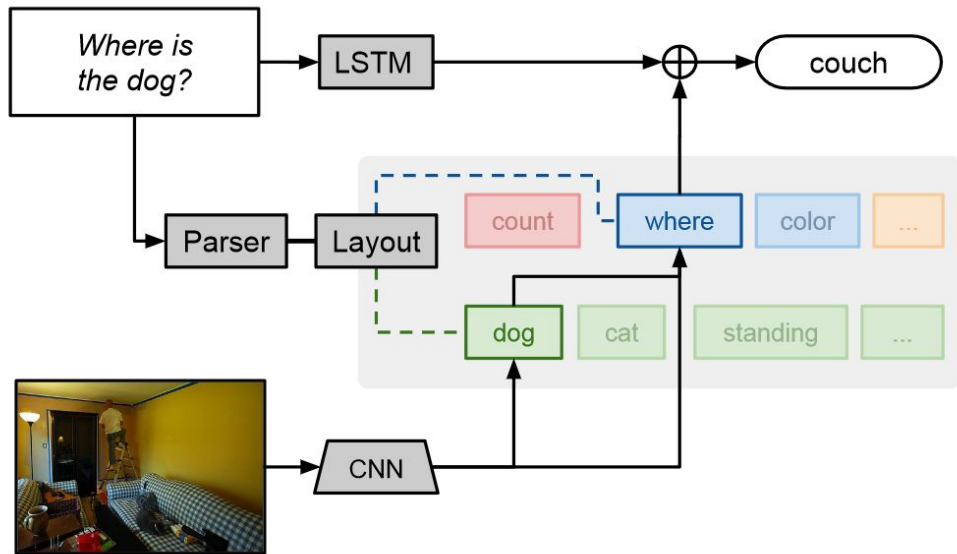
q is a natural language question, y is a label (or boolean), w can be visual **or semantic**.

Would like to work well with a small amount of data, but still benefit from significant amounts of data.

Neural Module Networks

Answer a question over an input (image only), in two steps:

1. *Layout* a network from the question.
2. *Evaluate* the network on the input.





Neural Module Networks

Two large weaknesses:

1. What if we don't have an image as input?
2. What if dependency parsing results in a bad network layout?

What if we don't have an image as input?



Replace Image with “World”

- The “World” is an arbitrary set of vectors.
- Still use attention across the vectors.
- Treat image as world by operating after the CNN.
- NMN modules assume CNN / Image!



New Modules!

Neural Module Network	Dynamic Neural Module Network
$\text{attend}[\text{word}]: \text{Image} \rightarrow \text{Attention}$	$\text{find}[\text{word}]: (\text{World}) \rightarrow \text{Attention}$
	$\text{lookup}[\text{word}]: () \rightarrow \text{Attention}$
$\text{re-attend}[\text{word}]: \text{Attention} \rightarrow \text{Attention}$	$\text{relate}[\text{word}]: (\text{World}) \text{ Attention} \rightarrow \text{Attention}$
$\text{combine}[\text{word}]: \text{Attention} \times \text{Att.} \rightarrow \text{Attention}$	$\text{and}: \text{Attention}^* \rightarrow \text{Attention}$
$\text{classify}[\text{word}]: \text{Image} \times \text{Attention} \rightarrow \text{Label}$	$\text{describe}[\text{word}]: (\text{World}) \text{ Attention} \rightarrow \text{Labels}$
$\text{measure}[\text{word}]: \text{Attention} \rightarrow \text{Label}$	$\text{exists}: \text{Attention} \rightarrow \text{Labels}$



Attend → Find

Neural Module Network	Dynamic Neural Module Network
$\text{attend}[\text{word}]: \textit{Image} \rightarrow \textit{Attention}$	$\text{find}[\text{word}]: (\textit{World}) \rightarrow \textit{Attention}$
A convolution.	“An MLP:” $\text{softmax}(a \odot \sigma(Bv_i \oplus CW \oplus d))$
$\text{attend}[\text{dog}]$	$\text{find}[\text{dog}]$ or $\text{find}[\text{city}]$
Generates an attention over the <i>Image</i> .	Generates an attention over the <i>World</i> .



“ “ → Lookup

Neural Module Network	Dynamic Neural Module Network
	lookup[word]: () → <i>Attention</i>
	A know relation: $e^{f(i)}$
	lookup[Georgia]
	For words with constant attention vectors.



Re-attend → Relate

Neural Module Network	Dynamic Neural Module Network
re-attend[word]: <i>Attention</i> → <i>Attention</i>	relate[word]: (<i>World</i>) <i>Attention</i> → <i>Attention</i>
(FC → ReLU) x 2	$\text{softmax}(a \odot \sigma(Bv_i \oplus CW \oplus Dw(h) \oplus e))$
re-attend[above]	relate[above] OR relate[in]
Generates a new attention over the <i>Image</i> .	Generates a new attention over the <i>World</i> .



Combine \rightarrow And

Neural Module Network	Dynamic Neural Module Network
combine[word]: $Attention \times Att. \rightarrow Attention$	and: $Attention^* \rightarrow Attention$
Stack \rightarrow Conv. \rightarrow ReLU	$h1 \circ h2 \circ \dots$
combine[except]	and
Combines two <i>Attentions</i> in an arbitrary way.	Multiplies attentions (analogous to set intersection).



Classify → Describe

Neural Module Network	Dynamic Neural Module Network
classify[word]: <i>Image</i> x <i>Attention</i> → <i>Label</i>	describe[word]: (<i>World</i>) <i>Attention</i> → <i>Labels</i>
Attend → FC → Softmax	$\text{softmax}(A\sigma(Bw(h) + v_i))$
classify[where]	describe[color] OR describe[where]
Transforms an <i>Image</i> and <i>Attention</i> into a <i>Label</i> .	Transforms a <i>World</i> and <i>Attention</i> into a <i>Label</i> .



Measure → Exists

Neural Module Network	Dynamic Neural Module Network
measure[word]: <i>Attention</i> → <i>Label</i>	exists: <i>Attention</i> → <i>Labels</i>
FC → ReLU → FC → Softmax	softmax((argmax h) a + b)
measure[exists]	exists
Transforms just an <i>Attention</i> into a <i>Label</i> .	Transforms just an <i>Attention</i> into a <i>Label</i> .

**What if dependency parsing
results in a bad network layout?**



New layout algorithm!

NMN

- Dependency parse
 - Leaf → attend
 - Internal (arity 1) → re-attend
 - Internal (arity 2) → combine
 - Root (yes/no) → measure
 - Root (other) → classify
- Layout of network strictly follows structure of dependency parse tree.

Dynamic-NMN

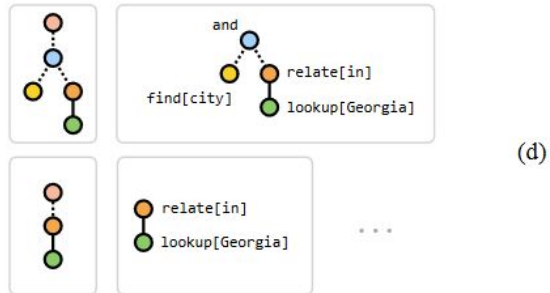
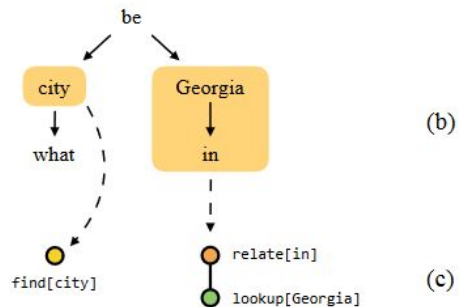
- Dependency parse
 - Proper nouns → lookup
 - Nouns & Verbs → find
 - Prepositional phrase → relate + find
- Generate candidate layouts from subsets of fragments.
 - and all fragments in subset
 - measure or combine
- “Rank” layouts with structure predictor.
- Use highly ranked layout.

New layout algorithm!

Only possible because “and” module has no parameters.

Structure predictor doesn't have any direct supervision. How can we train it?

What cities are in Georgia? (a)





Structure Predictor?

Computes $h_q(x)$ by passing LSTM over question.

Computes featurization $f(z_i)$ of i th layout.

Sample layout with probability $p(z_i | x; \theta_l) = \text{softmax}(a \cdot \sigma(B h_q(x) + C f(z_i) + d))$



How to train Structure Predictor?

Use a gradient estimate, as in REINFORCE (Williams, 1992).

Want to perform an SGD update with $\nabla J(\theta_l)$.

Estimate $\nabla J(\theta_l) = \mathbb{E}[\nabla \log p(z | x; \theta_l) \cdot r]$

Use reward $r = \log p(y | z, w; \theta_e)$

Step in direction $\nabla \log p(z | x; \theta_l) \cdot \log p(y | z, w; \theta_e)$

With small enough learning rate, estimate should converge.



New Dataset: GeoQA (+ Q)

- Entirely semantic: database of relations.
- Very small: 263 examples.
- (+ Q) adds quantification questions (e.g. What cities are in Texas? → Are there any cities in Texas?)
- State of the art results.
 - Compared to 2013 baseline and NMN.

Model	Accuracy	
	GeoQA	GeoQA+Q
LSP-F	48	–
LSP-W	51	–
NMN	51.7	35.7
D-NMN	54.3	42.9

Table 2: Results on the GeoQA dataset, and the GeoQA dataset with quantification. Our approach outperforms both a purely logical model (LSP-F) and a model with learned perceptual predicates (LSP-W) on the original dataset, and a fixed-structure NMN under both evaluation conditions.



Old Dataset: VQA

- Need to add “passthrough” to final hidden layer.
- Once again uses pre-trained VGG network.
- Slightly improved state of the art.

	test-dev				test-std
	Yes/No	Number	Other	All	All
Zhou (2015)	76.6	35.0	42.6	55.7	55.9
Noh (2015)	80.7	37.2	41.7	57.2	57.4
Yang (2015)	79.3	36.6	46.1	58.7	58.9
NMN	81.2	38.0	44.0	58.6	58.7
D-NMN	81.1	38.6	45.5	59.4	59.4

Table 1: Results on the VQA test server. NMN is the parameter-tying model from Andreas et al. (2015), and D-NMN is the model described in this paper.



Weaknesses?

- Can only generate very flat layouts, with only one conjunction or quantifier.
- Gradient estimate probably much more expensive / unstable than true gradient.
- Not any simpler than NMN, which are already considered complex.
- Similar in spirit but not implementation to Neural Symbolic VQA (Yi et. al. 2018).
- Much more complex than Relation Networks (Santoro et. al. 2017).

Questions? Discussion.

