Diagnostic classification and symbolic guidance to understand and improve recurrent neural networks



Dieuwke Hupkes and Willem Zuidema ILLC, University of Amsterdam



The task To understand how **recurrent networks** can implement and learn deeply hierarchical **compositional structures**, we train them on a toy task that involves computing the meaning of expressions with a hierarchical compositional semantics.

Network performance

GRU average

GRU best

LSTM best

--- LSTM average

Arithmetic language

-two minus ((two plus eight) minus (six plus one))

Sentences consist of digits and operators. Brackets indicate compositional structure.







Symbolic hypotheses																																	
<pre>minus_scope3+ minus_scope2+ minus_scope1+ close_minus_scope1+</pre>	0	0	0	0	1 1	1 1	1 1	1 1 2	1 1 3	1 1 3	1 1 3	1 1 4	1 1 4	1 1 1 4	1 1 1 4	1 1 3	1 1 2	1 1 2	1 1 1 3	1 1 1 3	1 1 1 3	1 1 1 3	1 1 2	1 1 1	1 0	0	0	1 1	1 1	1 1 1	1 1 1	1 0	0
	((-2	-	(6	-	((8	+	(-3	-	10))	-	(-2	-	10))))	-	(1	-	-8))
mode switch_mode	+	+	+	- 1	_	_	+ 1	+	+	+	+	+	+	- 1	_	+ 1	+	- 1	_	_	+ 1	+	- 1	+ 1	- 1	+ 1	- 1	_	_	+ 1	+	- 1	+ 1

$(\circ\circ)$

Symbolic Guidance

0

 $\left(\circ \circ \circ \right)$

Using the findings of *diagnostic classification* to provide additional symbolically motivated

