

Revisit Systematic Generalization via Meaningful Learning

Ning Shi, Boxin Wang, Wei Wang, Xiangyu Liu, Zhouhan Lin

ning.shi@ualberta.ca, boxinw2@illinois.edu, {luyang.w.w, eason.lxy}@alibaba-inc.com, lin.zhouhan@gmail.com

Introduction

Humans can systematically generalize to novel compositions of existing concepts. Recent studies argue that neural networks appear inherently ineffective in such cognitive capacity, leading to a *pessimistic* view and a lack of attention to *optimistic* results.

In contrast, the successful one-shot generalization in the turn-left experiment on the Simplified CommAI Navigation (SCAN) task reveals the potential of seq2seq recurrent networks in controlled environments (Lake and Baroni, 2018).

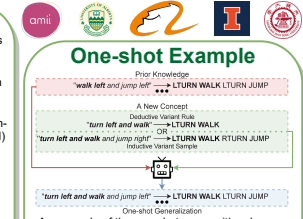
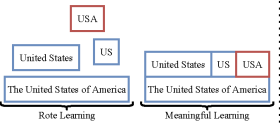
Question by Lake and Baroni (2018) on page 8:

"What are, precisely, the generalization mechanisms that subvert the networks' success in these experiments?"

Meaningful Learning

In educational psychology, *meaningful learning* refers to learning new concepts by relating them to old ones (Ausubel, 1963).

On the contrary, *rote learning* stands for learning new concepts without the consideration of relationships.



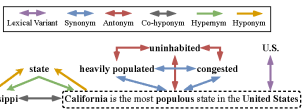
An example of the one-shot compositional generalization from the old concept "walk left" to the new one "turn left and walk" in SCAN.

From - "walk left and jump left"
To - "turn left and walk and jump left"

Old concept - "walk left"
New concept - "turn left and walk"

Connect both by pointing to **LTURN WALK**:

- Inductive Learning
- Deductive Learning



Inductive Learning

Inductive learning is a *bottom-up* approach from the more specific to the more general. In grammar teaching, inductive learning is a rule-discovery approach starting with the presentation of specific examples from which a general rule can be inferred.

Deductive Learning

Deductive Learning, the opposite of inductive learning, is a *top-down* approach from the more general to the more specific. As a rule-driven approach, teaching in a deductive manner often begins with presenting a general rule followed by specific examples in practice where the rule is applied

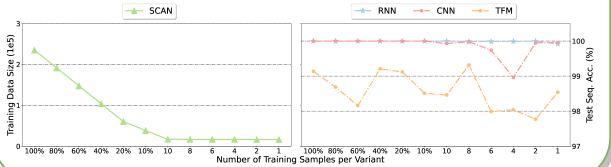
Data	Primitive	Variant	#Variants	Prompt
SCAN	jump	jump_0	10	[concept] twice
GEO	new york city mississippi river at dinner	houston city red river kansas tulsa	39 9 49 8	how many people in [concept] how long is [concept] where is [concept] what states capital is [concept]
ADV	a history of american film aaron swartz caspis 100	advanced ai techniques cargo aaronsw 171	5424 5492 51720 51895	who teaches [concept] ? does [concept] give upper-level courses ? name core courses for [concept] can undergrads take [concept] ?

Data	Primitive	Semantic Link	Variant	Primitive Rule	Concept Rule
SCAN	jump look run walk	Lexical Variant	jump_0 look_0 run_0 walk_0	jump → ELOP look → ELOP run → ELOP walk → ELOP	jump_0 → ELOP look_0 → ELOP run_0 → ELOP walk_0 → ELOP
GEO	new york city mississippi river at dinner	Co-hyponym	houston city red river kansas tulsa	new york city → CITY_NAME mississippi river → RIVER_NAME at → STATE_NAME dinner → CAPITAL_NAME	houston city → CITY_NAME red river → RIVER_NAME kansas → STATE_NAME tulsa → CAPITAL_NAME
ADV	a history of american film aaron swartz caspis 100	Co-hyponym	advanced ai techniques cargo aaronsw 171	a history of american film → TOPIC aaron swartz → INSTRUCTOR caspis → INSTRUMENT 100 → NUMBER	advanced ai techniques → TOPIC aaron swartz → INSTRUCTOR caspis → INSTRUMENT 171 → NUMBER

Systematic Generalization

Setup - we treat concepts in the initial data set as primitives and generate variant samples and rules accordingly. Next, we mix them up and construct a seq2seq task after a random split. We repeatedly train and evaluate models but slowly decrease the number of times they see each variant until one-shot learning.

Results - we observe there is *hardly* a performance drop for three representative model structures. This evidences that, with *semantic linking*, even canonical neural networks can generalize systematically to new concepts and compositions.



Semantic Linking Injection

We increase the difficulty of compositional generalization by excluding from the training set the primitive samples for inductive learning (left) and primitive rules for deductive learning (right).

Data	Model	Token Acc. %			Seq. Acc. %		
		Standard	Difficult	Challenging	Standard	Difficult	Challenging
SCAN	RNN	99.99 ± 0.03	99.89 ± 0.19	99.96 ± 0.02	99.95 ± 0.08	99.85 ± 0.08	99.80 ± 0.31
	CNN	99.96 ± 0.08	98.76 ± 0.54	98.89 ± 2.44	99.85 ± 0.34	99.32 ± 1.07	97.57 ± 5.24
	TFM	98.91 ± 0.78	98.90 ± 1.10	98.76 ± 0.85	97.35 ± 1.62	96.86 ± 2.64	96.38 ± 2.81
GEO	RNN	76.71 ± 8.42	75.69 ± 6.12	73.46 ± 3.05	44.96 ± 14.69	43.27 ± 13.47	36.77 ± 5.60
	CNN	87.99 ± 2.67	79.51 ± 6.03	77.40 ± 2.48	69.46 ± 5.78	51.20 ± 8.64	48.58 ± 3.40
	TFM	75.37 ± 7.84	75.11 ± 4.88	68.41 ± 4.76	45.93 ± 12.42	44.59 ± 9.76	36.93 ± 7.47
ADV	RNN	58.61 ± 6.18	59.74 ± 5.07	58.11 ± 5.82	36.18 ± 5.75	35.69 ± 6.05	35.45 ± 6.09
	CNN	57.83 ± 7.55	54.05 ± 5.74	53.66 ± 2.57	45.08 ± 9.32	42.14 ± 6.90	41.37 ± 4.04
	TFM	53.43 ± 2.80	51.51 ± 4.50	49.17 ± 2.58	42.59 ± 3.65	41.28 ± 4.35	38.88 ± 2.68

Data	Model	Token Acc. %		Seq. Acc. %	
		Standard	Difficult	Standard	Difficult
SCAN	RNN	99.48 ± 0.71	98.70 ± 0.92	98.27 ± 2.38	95.39 ± 2.72
	CNN	99.99 ± 0.01	98.59 ± 3.10	99.96 ± 0.03	96.66 ± 7.27
	TFM	96.90 ± 1.78	96.08 ± 2.21	91.94 ± 4.04	91.36 ± 5.80
GEO	RNN	54.44 ± 7.15	39.71 ± 18.38	13.61 ± 7.08	7.76 ± 5.34
	CNN	41.80 ± 3.38	41.07 ± 7.48	4.85 ± 4.66	4.04 ± 2.18
	TFM	67.02 ± 6.91	65.97 ± 5.17	36.38 ± 10.08	31.57 ± 7.42
ADV	RNN	36.50 ± 7.66	36.42 ± 7.39	12.84 ± 4.31	12.66 ± 5.19
	CNN	43.51 ± 11.31	35.34 ± 14.68	32.33 ± 12.93	23.58 ± 16.04
	TFM	56.82 ± 3.79	53.33 ± 3.85	47.43 ± 3.71	43.24 ± 5.14

Proof of Concept

Model	IMSLI14						IMSLI15					
	En-De		De-En		En-Fr		Fr-En		En-De		De-En	
	BLEU	SacreBLEU	BLEU	SacreBLEU	BLEU	SacreBLEU	BLEU	SacreBLEU	BLEU	SacreBLEU	BLEU	SacreBLEU
Baselines	26.98	24.88	30.18	32.62	38.06	42.93	57.34	59.36	26.98	24.88	30.18	32.62
LSTM (Liu et al., 2015)	28.95	28.85	35.24	37.60	41.82	46.41	40.45	42.61	28.95	28.85	35.24	37.60
Dynamic Cnn. (Wu et al., 2019)	27.38	27.28	33.33	35.54	40.41	45.32	58.41	61.02	27.38	27.28	33.33	35.54
Vocabulary Augmentation	25.59 ¹ ₂₁	25.36 ¹ ₂₂	30.99 ¹ ₂₁	33.63 ¹ ₂₂	38.32 ¹ ₂₀	43.38 ¹ ₂₁	57.77 ¹ ₂₀	59.83 ¹ ₂₁	25.59 ¹ ₂₁	25.36 ¹ ₂₂	30.99 ¹ ₂₁	33.63 ¹ ₂₂
Transformer (Vaswani et al., 2017)	29.48 ¹ ₂₁	29.29 ¹ ₂₂	35.72 ¹ ₂₁	38.07 ¹ ₂₂	43.19 ¹ ₂₁	48.68 ¹ ₂₁	61.86 ¹ ₂₀	63.13 ¹ ₂₁	29.48 ¹ ₂₁	29.29 ¹ ₂₂	35.72 ¹ ₂₁	38.07 ¹ ₂₂
Dynamic Cnn. (Wu et al., 2019)	27.68 ¹ ₂₁	27.50 ¹ ₂₂	33.62 ¹ ₂₁	36.09 ¹ ₂₂	40.87 ¹ ₂₀	45.99 ¹ ₂₁	59.85 ¹ ₂₀	61.86 ¹ ₂₁	27.68 ¹ ₂₁	27.50 ¹ ₂₂	33.62 ¹ ₂₁	36.09 ¹ ₂₂

Model	Geography				Advising			
	Token Acc. %	Seq. Acc. %	Token Acc. %	Seq. Acc. %	Token Acc. %	Seq. Acc. %	Token Acc. %	Seq. Acc. %
Baselines	89.08	17.39	69.81	9.68	92.22	3.84	40.41	6.11
CNN	96.45	79.74	78.44	55.01	98.74	81.62	81.74	51.13
TFM	99.45	84.95	80.24	49.92	99.68	76.90	78.51	28.67
Entity Segmentation	29.96	72.39 ¹ ₂₀	15.09 ¹ ₂₁	88.82	30.87	71.17 ¹ ₂₀	16.06 ¹ ₂₀	56.02 ¹ ₂₁
RNN	87.47	76.03	80.92 ¹ ₂₀	40.97 ¹ ₂₁	98.74	81.61	84.56 ¹ ₂₀	56.02 ¹ ₂₁
CNN	97.54	76.03	80.92 ¹ ₂₀	40.97 ¹ ₂₁	98.74	81.61	84.56 ¹ ₂₀	56.02 ¹ ₂₁
TFM	99.50	83.75	81.89 ¹ ₂₀	54.84 ¹ ₂₁	99.57	86.94	84.28 ¹ ₂₀	55.08 ¹ ₂₁

@BlackboxNLP

This work was supported by Shining Lab, Learnable, Inc., and Alibaba Group.