

# Hu et al., 2020 Sinha et al., 2019

---

Greta Tuckute & Kamoya K Ikhofua

MIT Fall 2020, 6.884  
*Symbolic Generalization*

# Motivation

Natural language understanding systems to generalize in a systematic and robust way

- Diagnostic tests - how can we probe these generalization abilities?
  - **Syntactic generalization** (Hu et al., 2020, “SG”) and **logical reasoning** (Sinha et al., 2019, “CLUTRR”)
- Evaluation metrics for language models?

# SG: Man shall not live by perplexity alone

Perplexity **is not sufficient** to check for human-like syntactic knowledge:

- It basically measures the probability of seeing some collection of words together
- However some words which are rarely seen together are grammatically correct
- *Colorless green ideas sleep furiously* (Chomsky, 1957)
- Need a **more fine-grained** way to assess learning outcomes of neural language models

# SG: Paradigm

Assess NL models on custom sentences designed using psycholinguistic and syntax literature/methodology

- Compare critical sentence regions NOT full-sentence probabilities.
- Factor out confounds (e.g token lexical frequency, n-gram statistics)

# SG: Paradigm

- Cover the scope of syntax phenomena: 16/47 (Carnie et al., 2012)
- Group syntax phenomena into 6 circuits based on processing algorithm

# SG: Circuits

1. Agreement
2. Licensing
3. Garden-Path Effects
4. Gross Syntactic Expectation
5. Center Embedding
6. Long-Distance Dependencies

## SG: Agreement

- (A) The farmer that the clerks embarrassed  $\text{knows}_{V_{sg}}$  many people.
- (B) \*The farmer that the clerks embarrassed  $\text{know}_{V_{pl}}$  many people.
- (C) The farmers that the clerk embarrassed  $\text{know}_{V_{pl}}$  many people.
- (D) \*The farmers that the clerk embarrassed  $\text{knows}_{V_{sg}}$  many people.

$$P_A(V_{sg}) > P_B(V_{pl}) \wedge P_C(V_{pl}) > P_D(V_{sg})$$

Chance is 25% (or up to 50%)

# SG: NPI Licensing

- The word “any” is a negative polarity item (NPI)
- The word “no” can license an NPI when it structurally commands it, such as in A

A) **No** managers that respected the guard have had **any** luck

>

B) \*The managers {that respected **no** guard} have had **any** luck

(Reflexive Pronoun Licensing was also included in sub-class suites)



# SG: NPI Licensing

(A) No managers that respected the guard have

<sup>NPI</sup>  
had any luck. [+NEG, -DISTRACTOR]

(B) \*The managers that respected no guard have

<sup>NPI</sup>  
had any luck. [-NEG, +DISTRACTOR]

(C) \*The managers that respected the guard have

<sup>NPI</sup>  
had any luck. [-NEG, -DISTRACTOR]

(D) No managers that respected no guard have

<sup>NPI</sup>  
had any luck. [+NEG, +DISTRACTOR]

$P_A(\text{NPI}) > P_C(\text{NPI}) \wedge P_D(\text{NPI}) > P_B(\text{NPI}) \wedge$

$P_A(\text{NPI}) > P_B(\text{NPI})$

Acceptable orderings:

ADBC

ADCB

DABC

DACB

ACDB (?)

Chance: 5/24

# SG: Reflexive Pronoun Licensing

- (A) The author that the senators liked hurt herself<sub>R<sub>sg.fem</sub></sub>.
- (B) \*The authors that the senator liked hurt herself<sub>R<sub>sg.fem</sub></sub>.
- (C) The authors that the senator liked hurt themselves<sub>R<sub>pl</sub></sub>.
- (D) \*The author~~X~~ that the senator liked hurt themselves<sub>R<sub>pl</sub></sub>.

$$P_A(R_{sg}) > P_B(R_{sg}) \wedge P_C(R_{pl}) > P_D(R_{pl})$$

Chance: 25%

# SG: NP/Z Garden-Paths

(A) !As the ship crossed the waters  $\overbrace{\text{remained}}^{V^*}$  blue and calm. [TRANS,NO COMMA]

(B) As the ship crossed, the waters  $\overbrace{\text{remained}}^{V^*}$  blue and calm. [TRANS,COMMA]

(C) As the ship drifted the waters  $\overbrace{\text{remained}}^{V^*}$  blue and calm. [INTRANS,NO COMMA]

(D) As the ship drifted, the waters  $\overbrace{\text{remained}}^{V^*}$  blue and calm. [INTRANS,COMMA]

$$S_A(V^*) > S_B(V^*) \wedge S_A(V^*) > S_C(V^*) \wedge S_A(V^*) - S_B(V^*) > S_C(V^*) - S_D(V^*)$$

# SG: Main-Verb Reduced Relative Garden-Paths

(A) !The child kicked in the chaos <sup>V\*</sup>found her way back home. [REDUCED, AMBIG]

(B) The child who was kicked in the chaos <sup>V\*</sup>found her way back home.





(C) The child forgotten in the chaos <sup>V\*</sup>found her way back home.

(D) The child who was forgotten in the chaos <sup>V\*</sup>found her way back home.

$$S_A(V^*) > S_B(V^*) \wedge S_A(V^*) > S_C(V^*) \wedge S_A(V^*) - S_B(V^*) > S_C(V^*) - S_D(V^*)$$

Chance is 25%

# SG: Gross Syntactic Expectation (*Subordination*)

- (A) The minister praised the building  . <sup>END</sup>
- (B) \*After the minister praised the building  . <sup>END</sup>
- (C) ??The minister praised the  
building,  it started to rain. <sup>MC</sup>
- (D) After the minister praised the  
building,  it started to rain. <sup>MC</sup>

$$P_A(\text{END}) > P_B(\text{END}) \wedge P_D(\text{MC}) < P_C(\text{MC})$$

## SG: Center Embedding

- (A) The painting<sub>N<sub>1</sub></sub> that the artist<sub>N<sub>2</sub></sub> who lived long ago painted<sub>V<sub>2</sub></sub> deteriorated<sub>V<sub>1</sub></sub>. [correct]
- (B) #The painting<sub>N<sub>1</sub></sub> that the artist<sub>N<sub>2</sub></sub> who lived long ago deteriorated<sub>V<sub>1</sub></sub> painted<sub>V<sub>2</sub></sub>. [incorrect]

$$P_A(V_2 V_1) > P_B(V_1 V_2)$$

$$P(\text{painted deteriorated} | \text{The painting that the artist}) >$$

$$P(\text{deteriorated painted} | \text{The painting that the artist})$$

# SG: Long Distance Dependencies

(A) I know that our uncle grabbed  $\overbrace{\text{the food}}^{\alpha}$  in front of the guests at the holiday party. [THAT, NO GAP]

(B) \*I know what our uncle grabbed  $\overbrace{\text{the food}}^{\alpha}$  in front of the guests at the holiday party. [WH, NO GAP]

(C) ??I know that our uncle grabbed  $\overbrace{\text{in front of}}^{\beta}$  the guests at the holiday party. [THAT, GAP]

(D) I know what our uncle grabbed  $\overbrace{\text{in front of}}^{\beta}$  in front of the guests at the holiday party. [WH, GAP]

$$S_B(\alpha) > S_A(\alpha) \wedge S_C(\beta) > S_D(\beta)$$

# SG: Pseudo-Clefting

(A) What the worker did was  $\overbrace{\text{board the plane.}}^{\text{VP}}$

(B) ?What the worker did was  $\overbrace{\text{the plane.}}^{\text{NP}}$

(C) What the worker repaired was  $\overbrace{\text{the plane.}}^{\text{NP}}$

(D) \*What the worker repaired was  
 $\overbrace{\text{board the plane.}}^{\text{VP}}$

$$S_D(\text{VP}) > S_A(\text{VP}) \wedge S_B(\text{NP}) > S_C(\text{NP})$$



# SG: Assessment

$\text{accuracy\_per\_test\_suite} = \text{correct predictions} / \text{total items}$

- Test for stability by including syntactically irrelevant but semantically plausible syntactic content before the critical region
  - E.g:
    - The keys to the cabinet on the left are on the table
    - \*The keys to the cabinet on the left is on the table
- Compare model class to dataset size

# SG: Score by Model Class

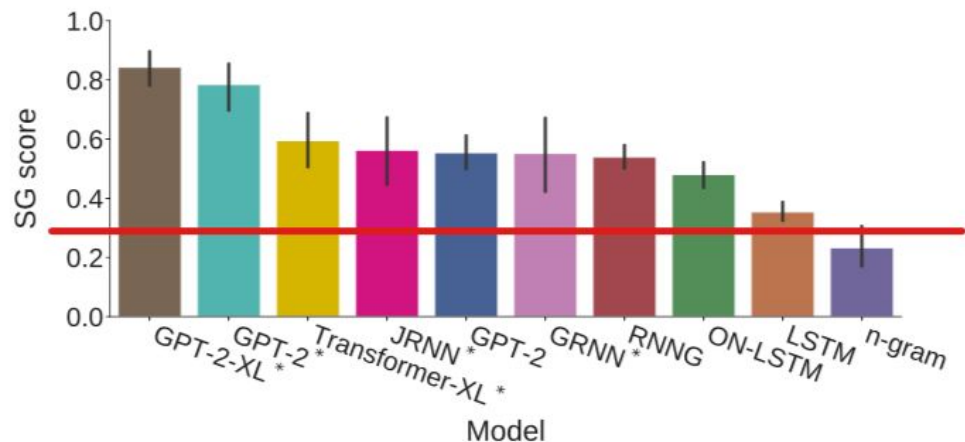
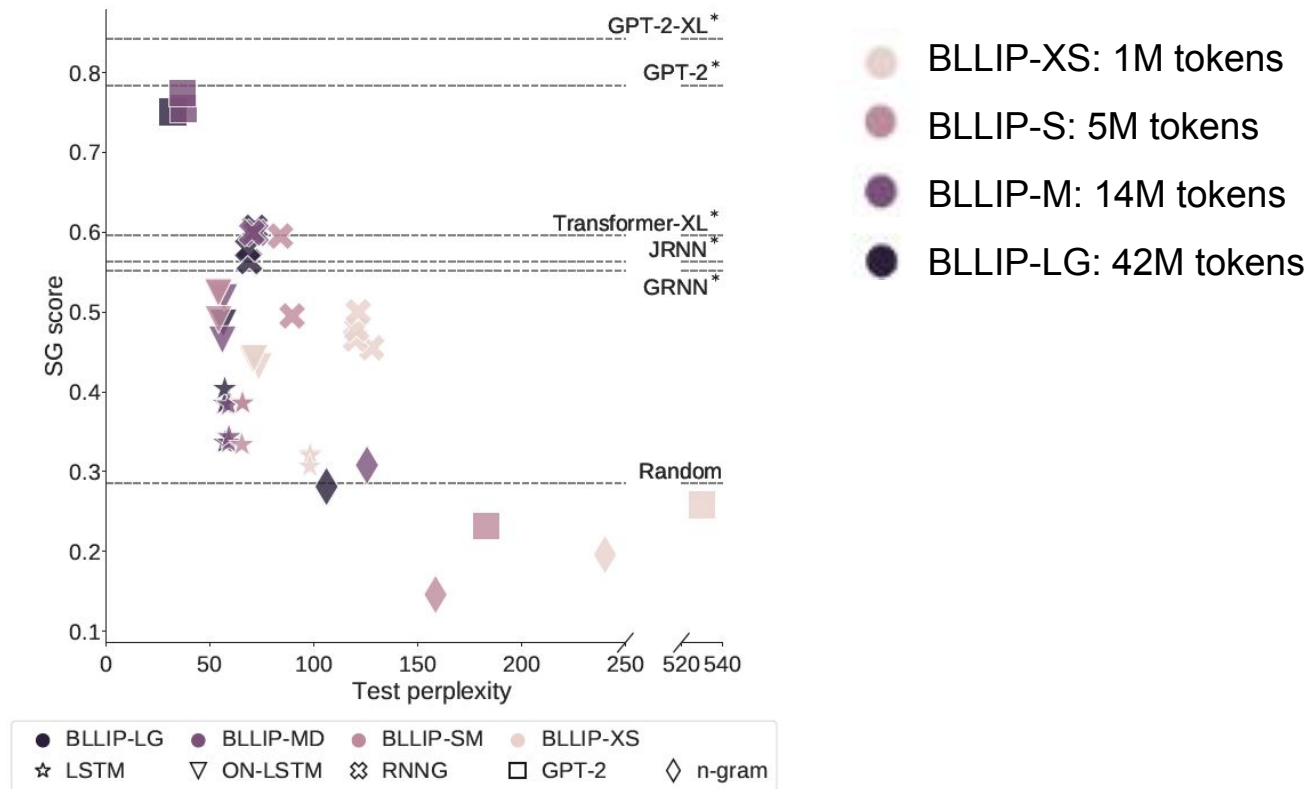
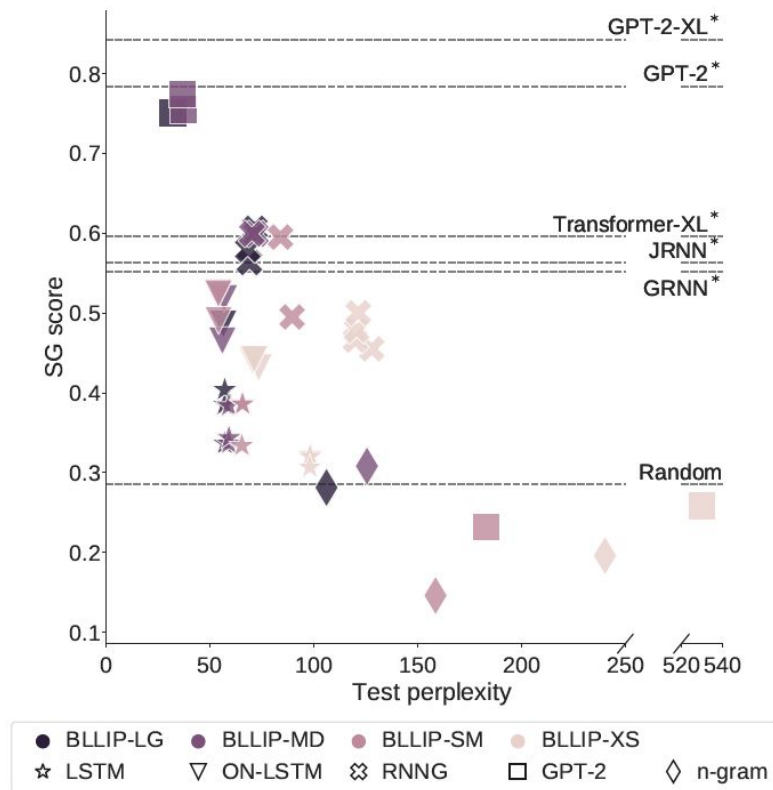


Figure 1: Average SG score by model class. Asterisks denote off-the-shelf models. Error bars denote bootstrapped 95% confidence intervals of the mean.

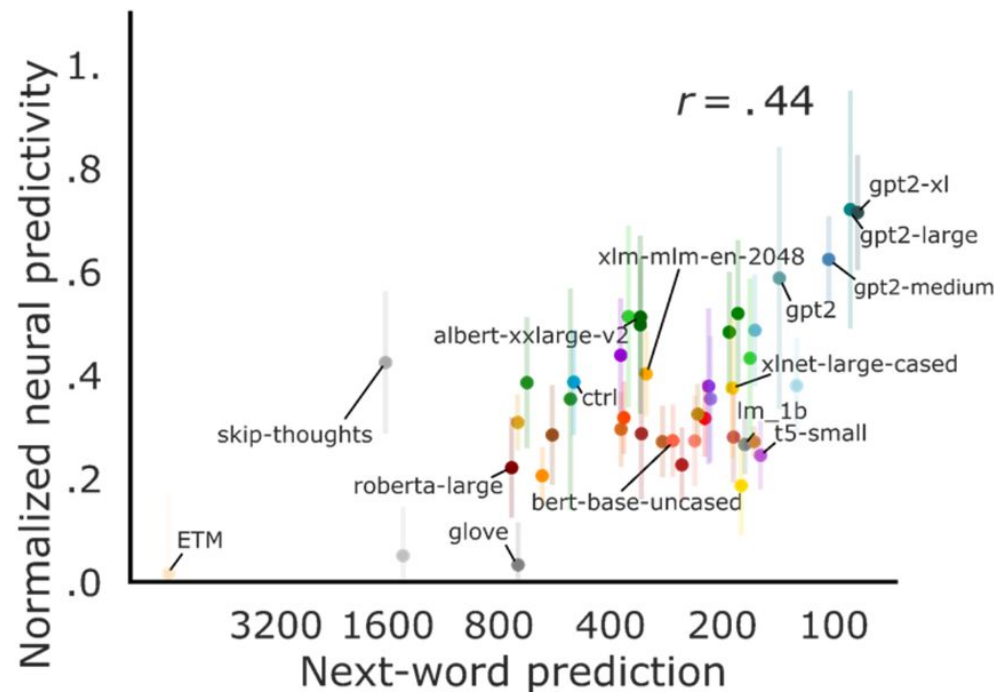
# SG: Perplexity and SG Score



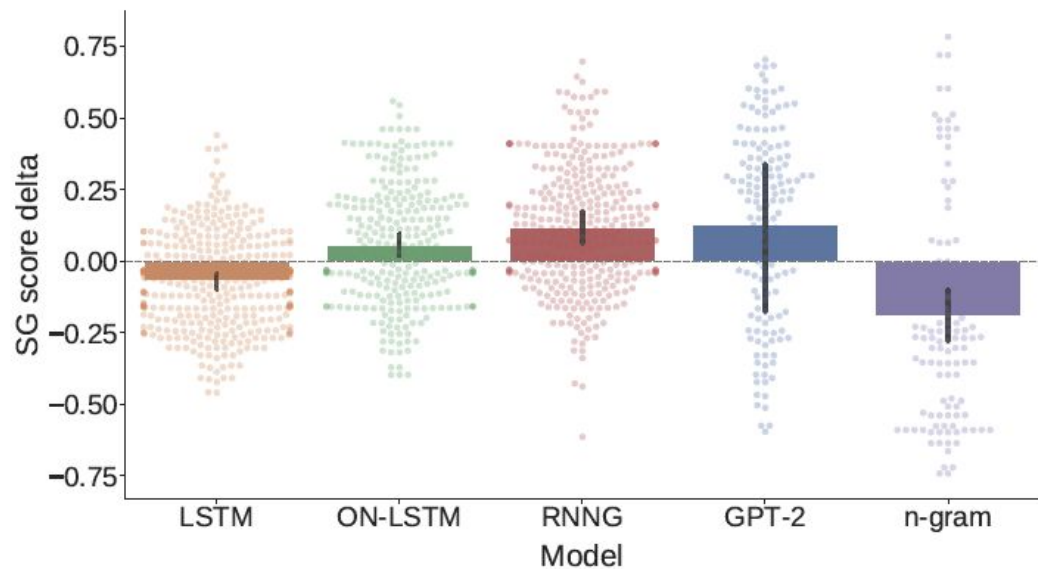
# SG: Perplexity and SG Score



# SG: Perplexity and Brain-Score

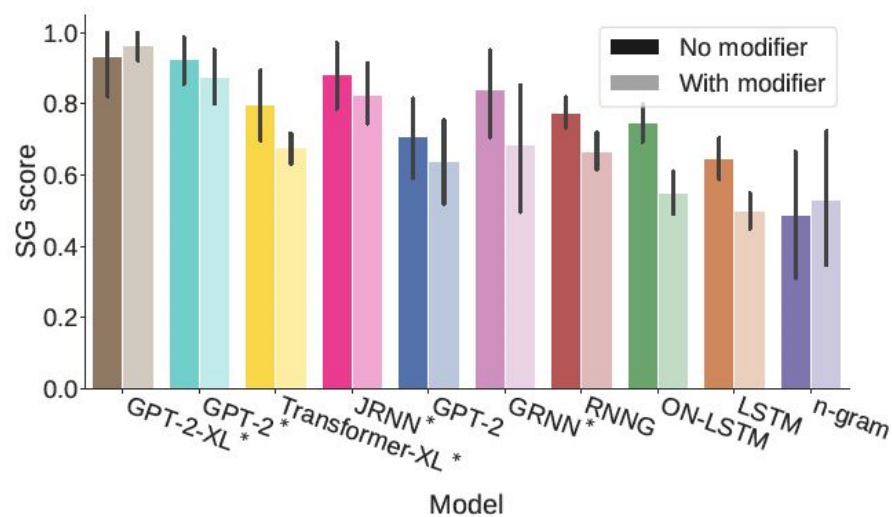


# SG: The Influence of Model Architecture

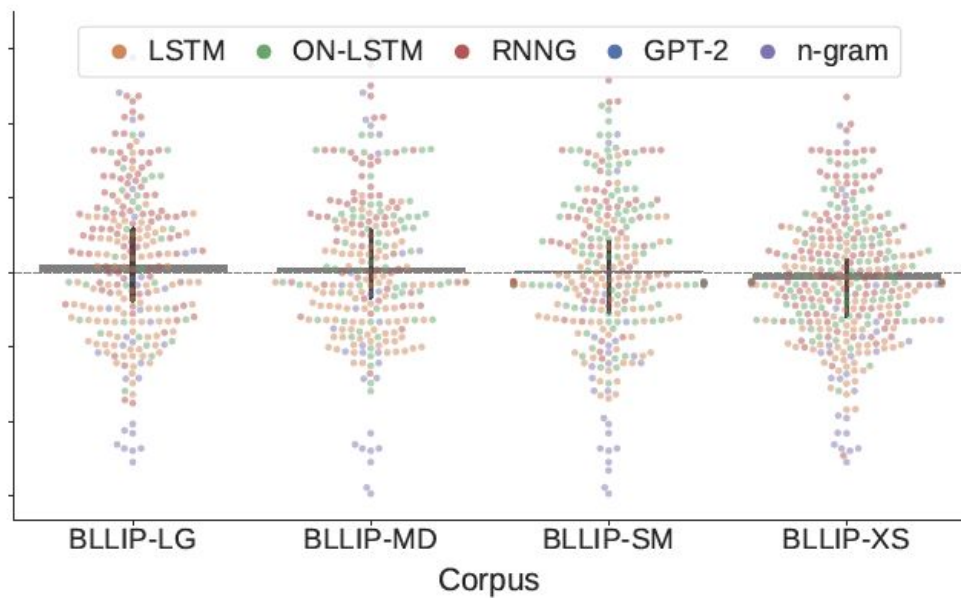


# SG: The Influence of Model Architecture

- Architectures as priors to the linguistic representation that can be developed
- Robustness depends on model architecture

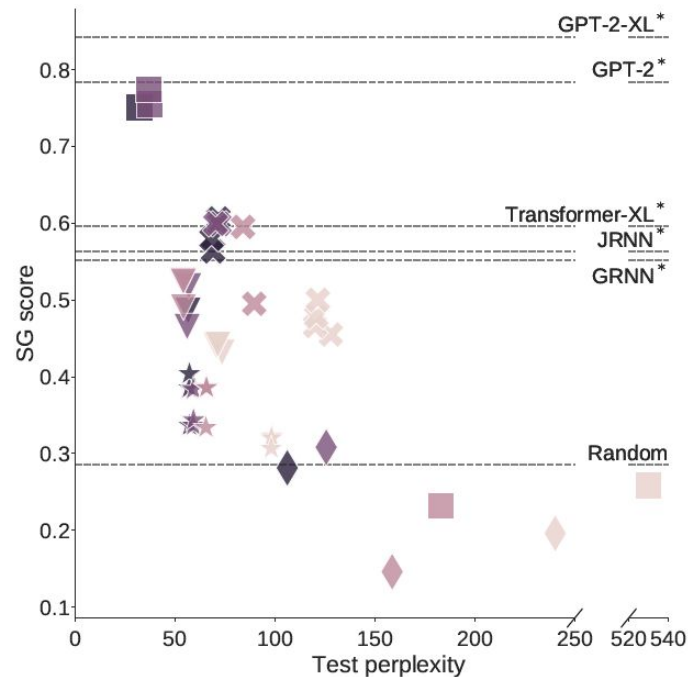
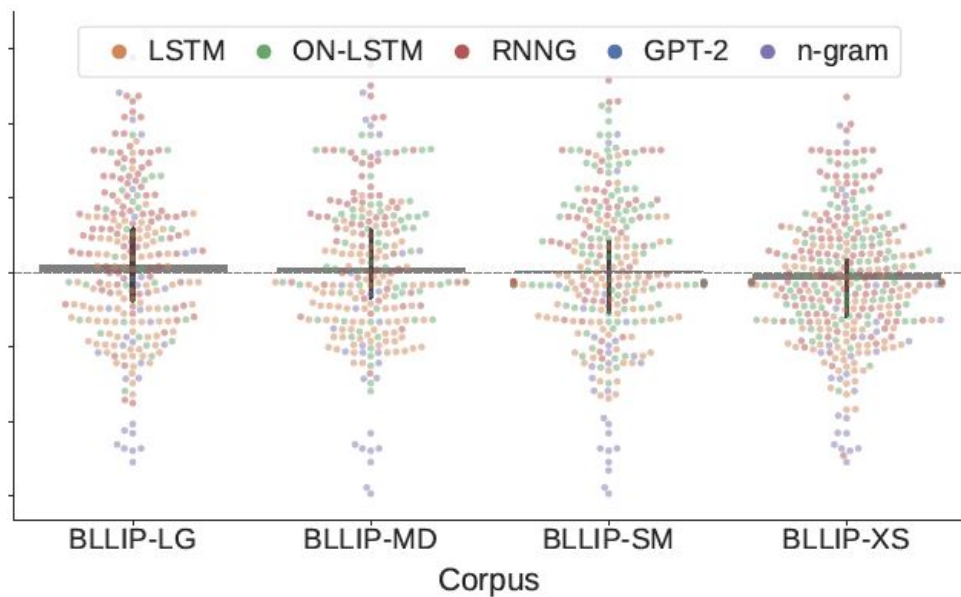


# SG: The Influence of Dataset Size





# SG: The Influence of Dataset Size



# SG: The Influence of Dataset Size

- Increasing amount of training data yields diminishing returns:
  - *“(...) require over 10 billion tokens to achieve human-like performance, and most would require trillions of tokens to achieve perfect accuracy – an impractically large amount of training data, especially for these relatively simple syntactic phenomena.”*  
(van Schijndel et al., 2019)
- Limited data efficiency
- Structured architectures or explicit syntactic supervision
- Humans? 11-27 million total words of input per year? (Hart & Risley, 1995; Brysbaert et al., 2016)

# SG: The Influence of Dataset Size

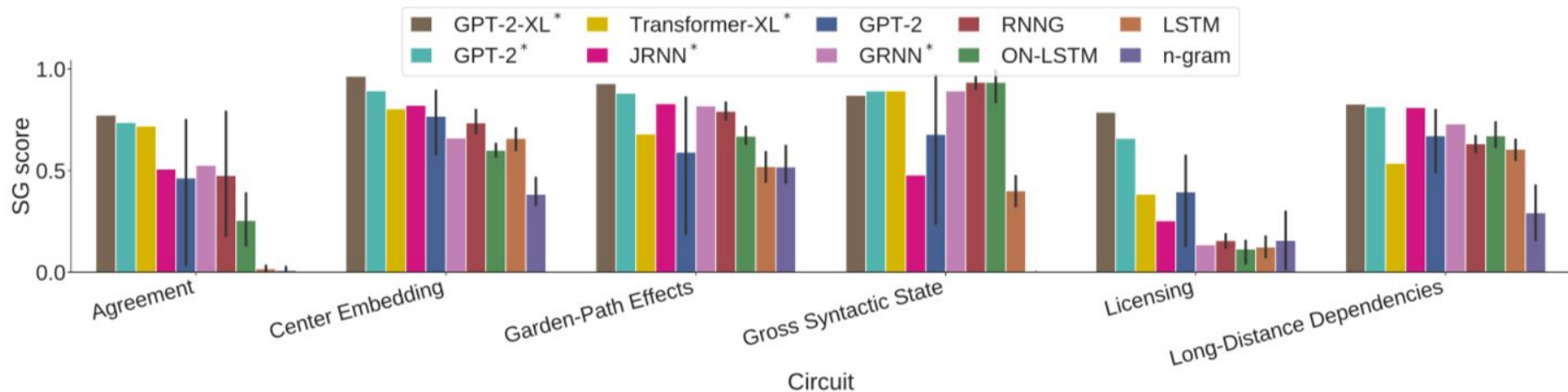


Figure 5: Evaluation results on all models, split across test suite circuits.

# CLUTRR: Motivation and Paradigm

- Compositional Language Understanding and Text-based Relational Reasoning
- Kinship inductive reasoning
- Unseen combinations of logical rules
- Model robustness

**Kristin** and her **son Justin** went to visit her **mother Carol** on a nice Sunday afternoon. They went out for a movie together and had a good time.



Q: How is **Carol** related to **Justin** ?

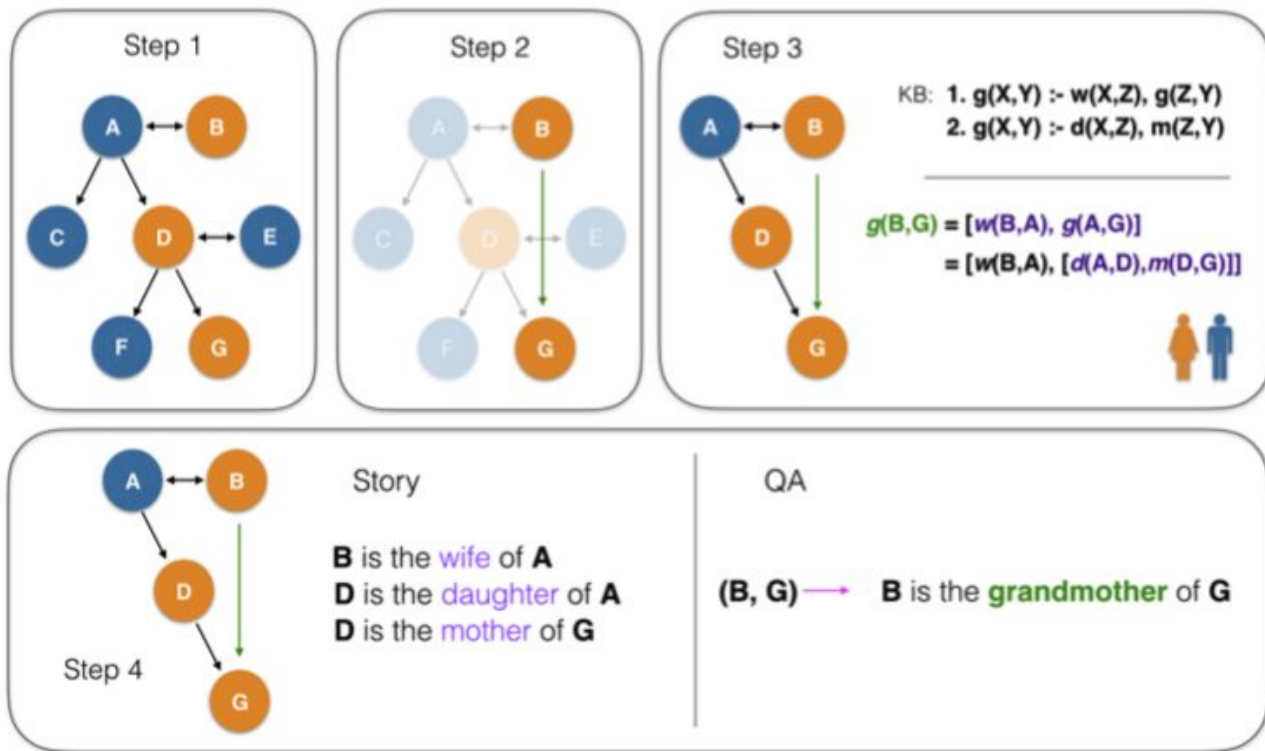
A: Carol is the **grandmother** of Justin



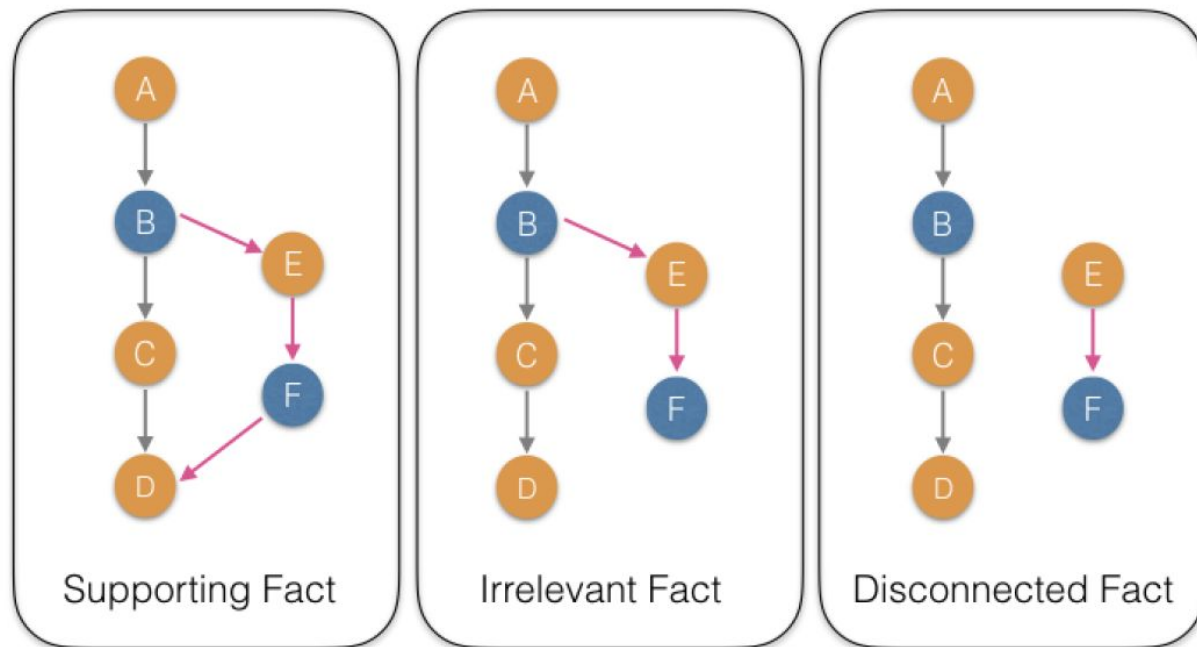
# CLUTRR: Motivation and Paradigm

- Productivity
  - $\text{mother}(\text{mother}(\text{mother}(\text{Justin}))) \sim$  great grandmother of Justin
- Systematicity
  - Only certain sets allowed with symmetries:  $\text{son}(\text{Justin}, \text{Kristin}) \sim \text{mother}(\text{Kristin}, \text{Justin})$
- Compositionality
  - $\text{son}(\text{Justin}, \text{Kristin})$  consists of components
- Memory (compression)
- Children are not exposed to systematic dataset

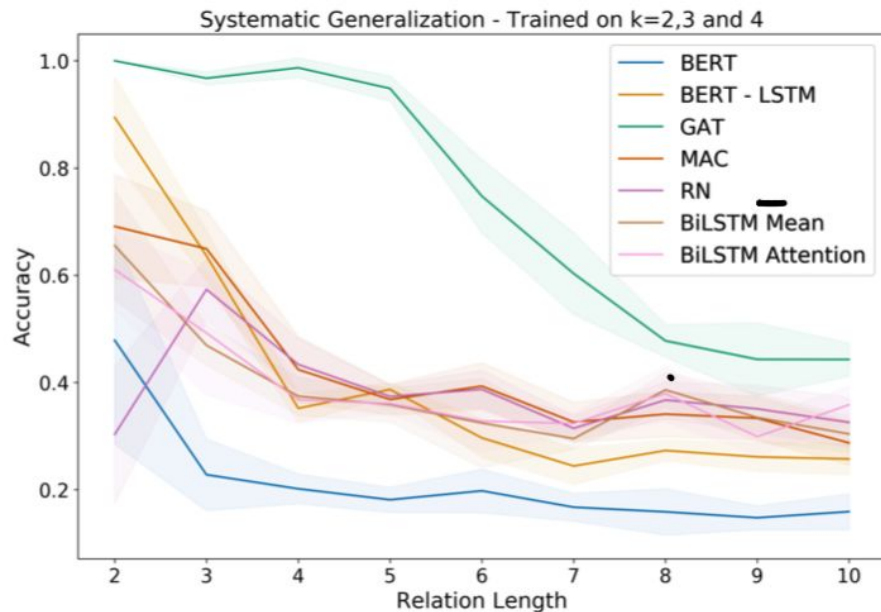
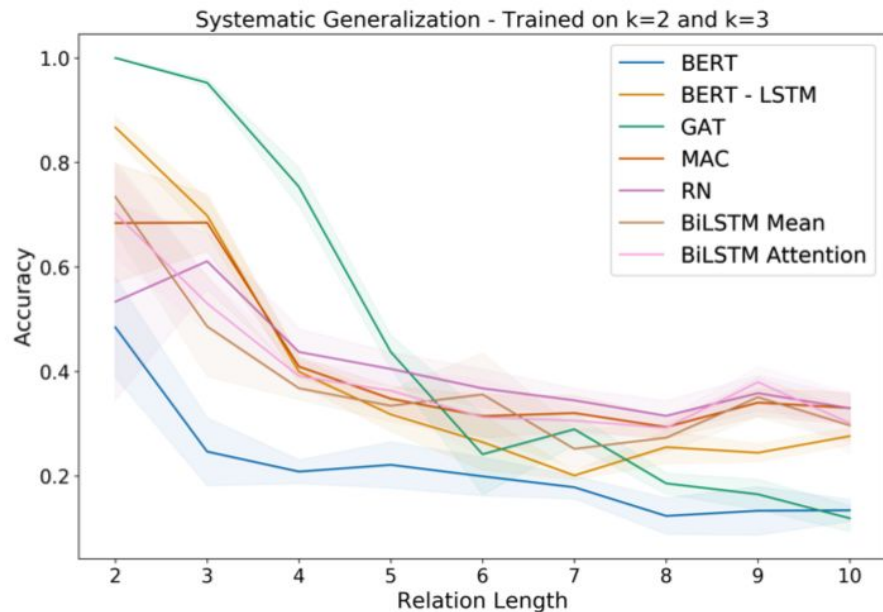
# CLUTRR: Dataset Generation & Paradigm



# CLUTRR: Model Robustness



# CLUTRR: Systematic Generalization





# CLUTRR: Model Robustness

Models		Unstructured models (no graph)						Structured model (with graph)
Training	Testing	BiLSTM - Attention	BiLSTM - Mean	RN	MAC	BERT	BERT-LSTM	GAT
Clean	Clean	0.58 $\pm$ 0.05	0.53 $\pm$ 0.05	0.49 $\pm$ 0.06	0.63 $\pm$ 0.08	0.37 $\pm$ 0.06	0.67 $\pm$ 0.03	<b>1.0</b> $\pm$ 0.0
	Supporting	<b>0.76</b> $\pm$ 0.02	0.64 $\pm$ 0.22	0.58 $\pm$ 0.06	0.71 $\pm$ 0.07	0.28 $\pm$ 0.1	0.66 $\pm$ 0.06	0.24 $\pm$ 0.2
	Irrelevant	0.7 $\pm$ 0.15	<b>0.76</b> $\pm$ 0.02	0.59 $\pm$ 0.06	0.69 $\pm$ 0.05	0.24 $\pm$ 0.08	0.55 $\pm$ 0.03	0.51 $\pm$ 0.15
	Disconnected	0.49 $\pm$ 0.05	0.45 $\pm$ 0.05	0.5 $\pm$ 0.06	0.59 $\pm$ 0.05	0.24 $\pm$ 0.08	0.5 $\pm$ 0.06	<b>0.8</b> $\pm$ 0.17
Supporting	Supporting	0.67 $\pm$ 0.06	0.66 $\pm$ 0.07	0.68 $\pm$ 0.05	0.65 $\pm$ 0.04	0.32 $\pm$ 0.09	0.57 $\pm$ 0.04	<b>0.98</b> $\pm$ 0.01
Irrelevant	Irrelevant	0.51 $\pm$ 0.06	0.52 $\pm$ 0.06	0.5 $\pm$ 0.04	0.56 $\pm$ 0.04	0.25 $\pm$ 0.06	0.53 $\pm$ 0.06	<b>0.93</b> $\pm$ 0.01
Disconnected	Disconnected	0.57 $\pm$ 0.07	0.57 $\pm$ 0.06	0.45 $\pm$ 0.11	0.4 $\pm$ 0.1	0.17 $\pm$ 0.05	0.47 $\pm$ 0.06	<b>0.96</b> $\pm$ 0.01
Average		<b>0.61</b> $\pm$ 0.08	0.59 $\pm$ 0.08	0.54 $\pm$ 0.07	<b>0.61</b> $\pm$ 0.06	0.30 $\pm$ 0.07	0.56 $\pm$ 0.05	<b>0.77</b> $\pm$ 0.09

# CLUTRR: Model Robustness (noisy training)

Models		Unstructured models (no graph)						Structured model (with graph)
Training	Testing	BiLSTM - Attention	BiLSTM - Mean	RN	MAC	BERT	BERT-LSTM	GAT
Supporting	Clean	0.38 $\pm$ 0.04	0.32 $\pm$ 0.04	0.4 $\pm$ 0.09	0.45 $\pm$ 0.03	0.19 $\pm$ 0.06	0.39 $\pm$ 0.06	<b>0.92</b> $\pm$ 0.17
	Supporting	0.67 $\pm$ 0.06	0.66 $\pm$ 0.07	0.68 $\pm$ 0.05	0.65 $\pm$ 0.04	0.32 $\pm$ 0.09	0.57 $\pm$ 0.04	<b>0.98</b> $\pm$ 0.01
	Irrelevant	0.44 $\pm$ 0.03	0.39 $\pm$ 0.03	<b>0.51</b> $\pm$ 0.08	0.46 $\pm$ 0.09	0.2 $\pm$ 0.06	0.36 $\pm$ 0.05	0.5 $\pm$ 0.23
	Disconnected	0.31 $\pm$ 0.21	0.25 $\pm$ 0.16	0.47 $\pm$ 0.08	0.41 $\pm$ 0.06	0.2 $\pm$ 0.08	0.32 $\pm$ 0.04	<b>0.92</b> $\pm$ 0.05
Irrelevant	Clean	0.57 $\pm$ 0.05	0.56 $\pm$ 0.05	0.46 $\pm$ 0.13	0.67 $\pm$ 0.05	0.24 $\pm$ 0.06	0.46 $\pm$ 0.08	<b>0.92</b> $\pm$ 0.0
	Supporting	0.38 $\pm$ 0.22	0.31 $\pm$ 0.16	0.61 $\pm$ 0.07	0.61 $\pm$ 0.04	0.27 $\pm$ 0.06	0.46 $\pm$ 0.04	<b>0.77</b> $\pm$ 0.12
	Irrelevant	0.51 $\pm$ 0.06	0.52 $\pm$ 0.06	0.5 $\pm$ 0.04	0.56 $\pm$ 0.04	0.25 $\pm$ 0.06	0.53 $\pm$ 0.06	<b>0.93</b> $\pm$ 0.01
	Disconnected	0.44 $\pm$ 0.26	0.54 $\pm$ 0.27	0.55 $\pm$ 0.05	0.61 $\pm$ 0.06	0.26 $\pm$ 0.03	0.45 $\pm$ 0.08	<b>0.85</b> $\pm$ 0.25
Disconnected	Clean	0.45 $\pm$ 0.02	0.47 $\pm$ 0.03	0.53 $\pm$ 0.09	0.5 $\pm$ 0.06	0.22 $\pm$ 0.09	0.44 $\pm$ 0.05	<b>0.75</b> $\pm$ 0.07
	Supporting	0.47 $\pm$ 0.03	0.46 $\pm$ 0.05	0.54 $\pm$ 0.03	0.58 $\pm$ 0.06	0.22 $\pm$ 0.06	0.38 $\pm$ 0.08	<b>0.78</b> $\pm$ 0.12
	Irrelevant	0.47 $\pm$ 0.05	0.48 $\pm$ 0.03	0.52 $\pm$ 0.04	0.51 $\pm$ 0.05	0.17 $\pm$ 0.04	0.38 $\pm$ 0.05	<b>0.56</b> $\pm$ 0.26
	Disconnected	0.57 $\pm$ 0.07	0.57 $\pm$ 0.06	0.45 $\pm$ 0.11	0.4 $\pm$ 0.1	0.17 $\pm$ 0.05	0.47 $\pm$ 0.06	<b>0.96</b> $\pm$ 0.01
Average		0.47 $\pm$ 0.08	0.46 $\pm$ 0.08	0.52 $\pm$ 0.07	<b>0.53</b> $\pm$ 0.06	0.23 $\pm$ 0.07	0.43 $\pm$ 0.05	<b>0.82</b> $\pm$ 0.09

Table 3: Testing the robustness of the various models when trained various types of noisy facts and evaluated on other noisy / clean facts. The types of noise facts (supporting, irrelevant and disconnected) are defined in Section 3.5 of the main paper.

# Future work & Perspectives

- Sub-word tokenization
- Active attention and reasoning
- Generalization across tasks
- Abstractions as probabilistic
- Architecture and dimensionality reduction

# References

Brysbaert, M., Stevens, M., Mander, P., & Keuleers, E. (2016). How Many Words Do We Know? Practical Estimates of Vocabulary Size Dependent on Word Definition, the Degree of Language Input and the Participant's Age. *Frontiers in psychology*, 7, 1116.

<https://doi.org/10.3389/fpsyg.2016.01116>

Hart, B., & Risley, T. R. (1995). *Meaningful differences in the everyday experience of young American children*. Baltimore, MD: Paul H. Brookes Publishing Company.

Schrimpf, M., Blank, I., Tuckute, G., Kauf, C., Hosseini, E. A., Kanwisher, N., Tenenbaum, J., Fedorenko, E (2020): Artificial Neural Networks Accurately Predict Language Processing in the Brain, *bioRxiv* 2020.06.26.174482; doi:

<https://doi.org/10.1101/2020.06.26.174482>.

Van Schijndel, M., Mueller, A., & Linzen, T. (2019). Quantity doesn't buy quality syntax with neural language models. *arXiv preprint arXiv:1909.00111*.

# Supplementary

# CLUTTR, Fig. 6

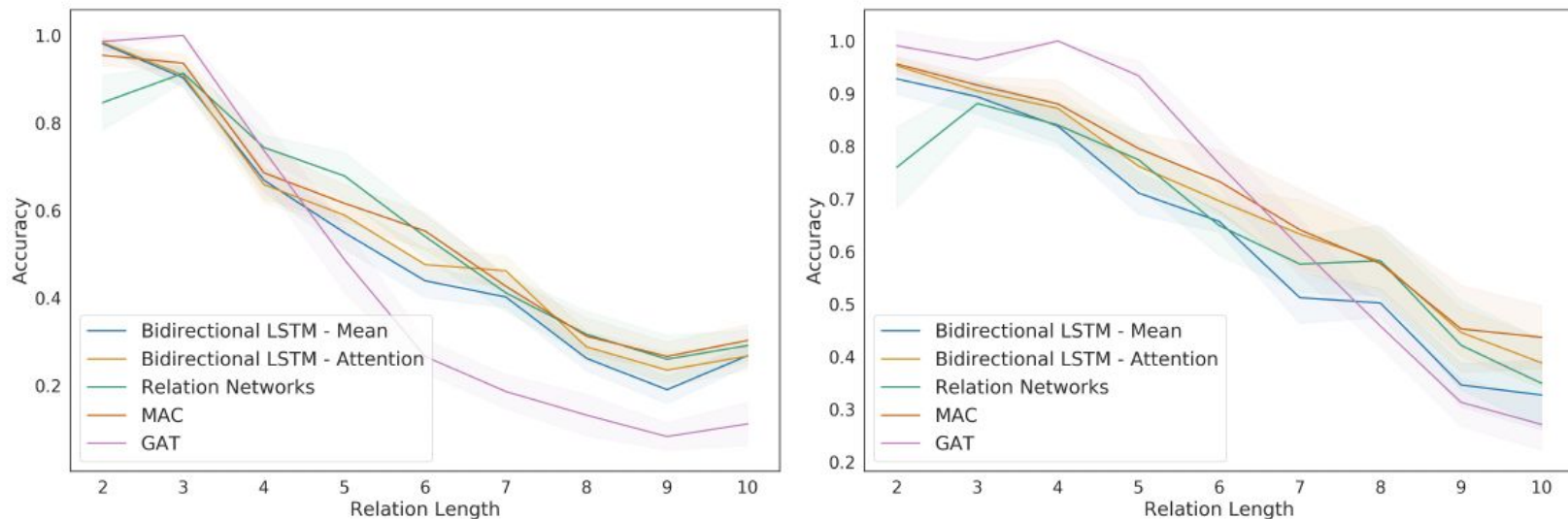


Figure 6: Systematic Generalizability of different models on CLUTTR-Gen task (having 20% less placeholders and without training and testing placeholder split), when **Left:** trained with  $k = 2$  and  $k = 3$  and **Right:** trained with  $k = 2, 3$  and  $4$

# CLUTTR, Table 5

Relation Length	Human Performance		Reported Difficulty
	Time Limited	Unlimited Time	
2	0.848	1	1.488 +- 1.25
3	0.773	1	2.41 +- 1.33
4	0.477	1	3.81 +- 1.46
5	0.424	1	3.78 +- 0.96
6	0.406	1	4.46 +- 0.87

Table 5: Human performance accuracies on CLUTRR dataset. Humans are provided the Clean-Generalization version of the dataset, and we test on two scenarios: when a human is given limited time to solve the task, and when a human is given unlimited time to solve the task. Regardless of time, our evaluators provide a score of difficulty of individual puzzles.

# CLUTTR, Table 4

Models		Unstructured models (no graph)						Structured model (with graph)
Training	Testing	BiLSTM - Attention	BiLSTM - Mean	RN	MAC	BERT	BERT-LSTM	GAT
Supporting	Clean	0.96 $\pm$ 0.01	<b>0.97</b> $\pm$ 0.01	0.88 $\pm$ 0.05	0.94 $\pm$ 0.02	0.48 $\pm$ 0.08	0.57 $\pm$ 0.08	0.92 $\pm$ 0.17
	Supporting	0.96 $\pm$ 0.03	0.96 $\pm$ 0.03	0.97 $\pm$ 0.01	0.97 $\pm$ 0.01	0.75 $\pm$ 0.07	0.88 $\pm$ 0.05	<b>0.98</b> $\pm$ 0.01
	Irrelevant	0.92 $\pm$ 0.02	<b>0.93</b> $\pm$ 0.01	0.9 $\pm$ 0.03	0.91 $\pm$ 0.01	0.56 $\pm$ 0.04	0.54 $\pm$ 0.06	0.5 $\pm$ 0.23
	Disconnected	0.8 $\pm$ 0.04	0.83 $\pm$ 0.04	0.76 $\pm$ 0.08	0.86 $\pm$ 0.04	0.27 $\pm$ 0.06	0.42 $\pm$ 0.08	<b>0.92</b> $\pm$ 0.05
Irrelevant	Clean	0.63 $\pm$ 0.02	0.61 $\pm$ 0.07	0.85 $\pm$ 0.09	0.8 $\pm$ 0.07	0.53 $\pm$ 0.09	0.44 $\pm$ 0.06	<b>0.92</b> $\pm$ 0.0
	Supporting	0.66 $\pm$ 0.03	0.64 $\pm$ 0.04	0.69 $\pm$ 0.06	0.76 $\pm$ 0.06	0.42 $\pm$ 0.08	0.43 $\pm$ 0.08	<b>0.77</b> $\pm$ 0.12
	Irrelevant	0.89 $\pm$ 0.04	0.86 $\pm$ 0.1	0.74 $\pm$ 0.11	0.78 $\pm$ 0.06	0.61 $\pm$ 0.1	0.83 $\pm$ 0.06	<b>0.93</b> $\pm$ 0.01
	Disconnected	0.64 $\pm$ 0.02	0.62 $\pm$ 0.05	0.72 $\pm$ 0.05	0.73 $\pm$ 0.04	0.41 $\pm$ 0.04	0.61 $\pm$ 0.05	<b>0.85</b> $\pm$ 0.25
Disconnected	Clean	0.9 $\pm$ 0.05	0.82 $\pm$ 0.12	<b>0.94</b> $\pm$ 0.02	0.93 $\pm$ 0.04	0.68 $\pm$ 0.07	0.64 $\pm$ 0.02	0.75 $\pm$ 0.07
	Supporting	0.87 $\pm$ 0.04	0.82 $\pm$ 0.05	0.85 $\pm$ 0.03	<b>0.88</b> $\pm$ 0.04	0.54 $\pm$ 0.08	0.5 $\pm$ 0.05	0.78 $\pm$ 0.12
	Irrelevant	<b>0.87</b> $\pm$ 0.03	0.85 $\pm$ 0.03	0.83 $\pm$ 0.03	0.87 $\pm$ 0.02	0.59 $\pm$ 0.09	0.58 $\pm$ 0.09	0.56 $\pm$ 0.26
	Disconnected	0.91 $\pm$ 0.04	0.91 $\pm$ 0.03	0.8 $\pm$ 0.17	0.71 $\pm$ 0.11	0.49 $\pm$ 0.1	0.79 $\pm$ 0.1	<b>0.96</b> $\pm$ 0.01
Average		0.83 $\pm$ 0.08	0.82 $\pm$ 0.08	0.83 $\pm$ 0.07	<b>0.84</b> $\pm$ 0.06	0.58 $\pm$ 0.07	0.60 $\pm$ 0.05	<b>0.82</b> $\pm$ 0.09

Table 4: Testing the robustness on toy placeholders of the various models when trained various types of noisy facts and evaluated on other noisy / clean facts. The types of noise facts (supporting, irrelevant and disconnected) are defined in Section 3.5 of the main paper.



# CLUTTR, Fig. 7

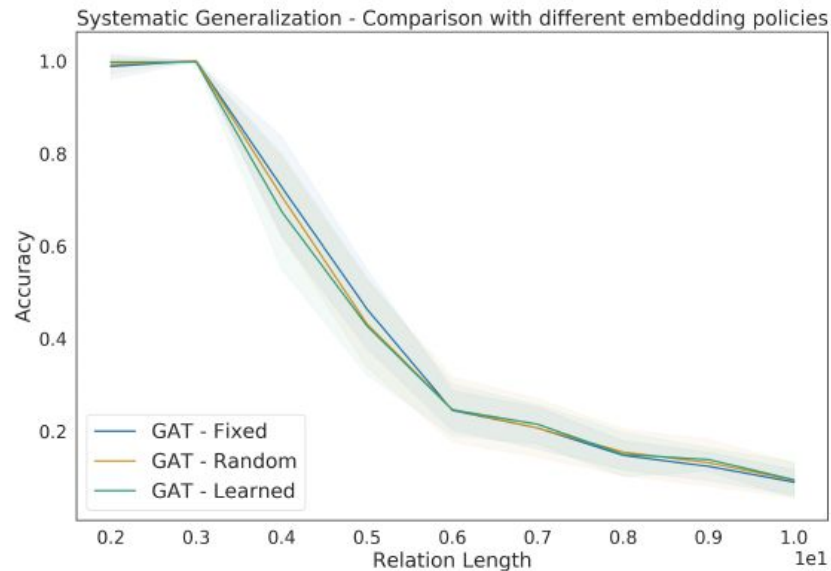
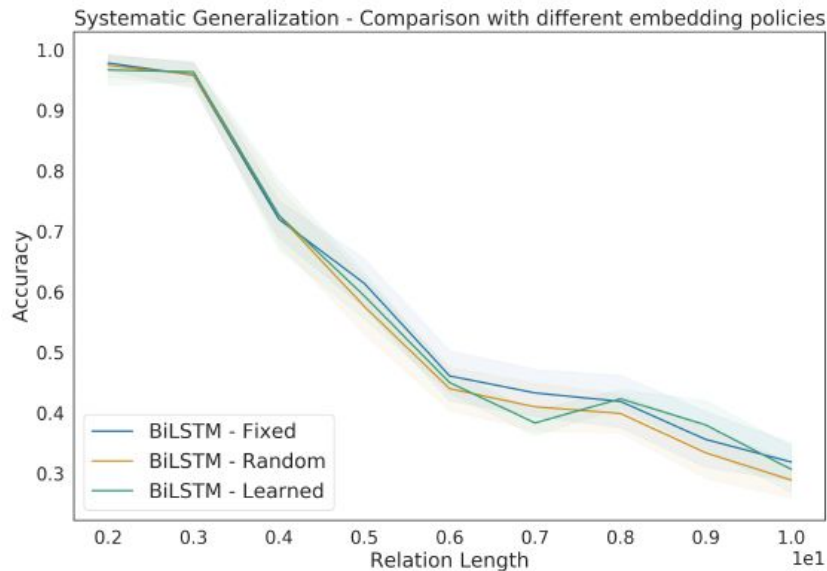


Figure 7: Systematic Generalization comparison with different Embedding policies

# Van Schijndel et al., 2019

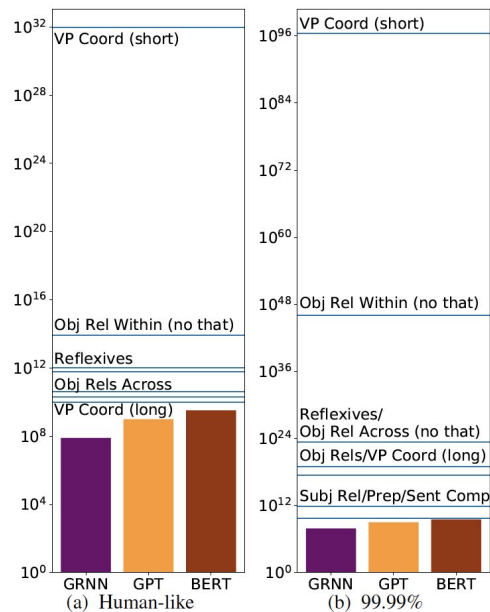


Figure 2: Lines depict number of training tokens needed for LSTMs to achieve human-like (left) or 99.99% accuracy (right) in each syntactic agreement condition, according to our estimates. Bars depict the amount of data on which each model was trained.