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

Greta Tuckute & Kamoya K Ikhofua

MIT Fall 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 test human-level understanding of syntax

# 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 knows$_{V_{sg}}$ many people.

(B) *The farmer that the clerks embarrassed know$_{V_{pl}}$ many people.

(C) The farmers that the clerk embarrassed know$_{V_{pl}}$ many people.

(D) *The farmers that the clerk embarrassed knows$_{V_{sg}}$ many people.

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

# 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
$$\underbrace{\text{had any}}_{\text{NPI}}\text{ luck. } [\text{+NEG,}-\text{DISTRACTOR}]$$

**(B)** *The managers that respected no guard have
$$\underbrace{\text{had any}}_{\text{NPI}}\text{ luck. } [-\text{NEG,+DISTRACTOR}]$$

**(C)** *The managers that respected the guard have
$$\underbrace{\text{had any}}_{\text{NPI}}\text{ luck. } [-\text{NEG,}-\text{DISTRACTOR}]$$

**(D)** No managers that respected no guard have
$$\underbrace{\text{had any}}_{\text{NPI}}\text{ luck. } [\text{+NEG,+DISTRACTOR}]$$

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

Acceptable orderings:

ADBC
ADCB
DABC
DACB
ACDB (?)

Chance: 5/24

# SG: NP/Z Garden-Paths



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

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

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

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

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

(Main Verb / Reduced Relative Clause paths were also included in sub-class suites)

# SG: Gross Syntactic Expectation

(A) The minister praised the building . <span>END</span>

(B) *After the minister praised the building . <span>END</span>

(C) ??The minister praised the building, it started to rain. <span>MC</span>

(D) After the minster praised the building, it started to rain. <span>MC</span>

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

# SG: Center Embedding

The paintings that the artist painted deteriorated

>

*The paintings that the artist deteriorated painted

# SG: Long Distance Dependencies

The **keys** to the cabinet **are** on the table

**>**

*The **keys** to the cabinet **is** on the table

# SG: Cleft

The **keys** to the cabinet **are** on the table

**>**

*The **keys** to the cabinet **is** on the table
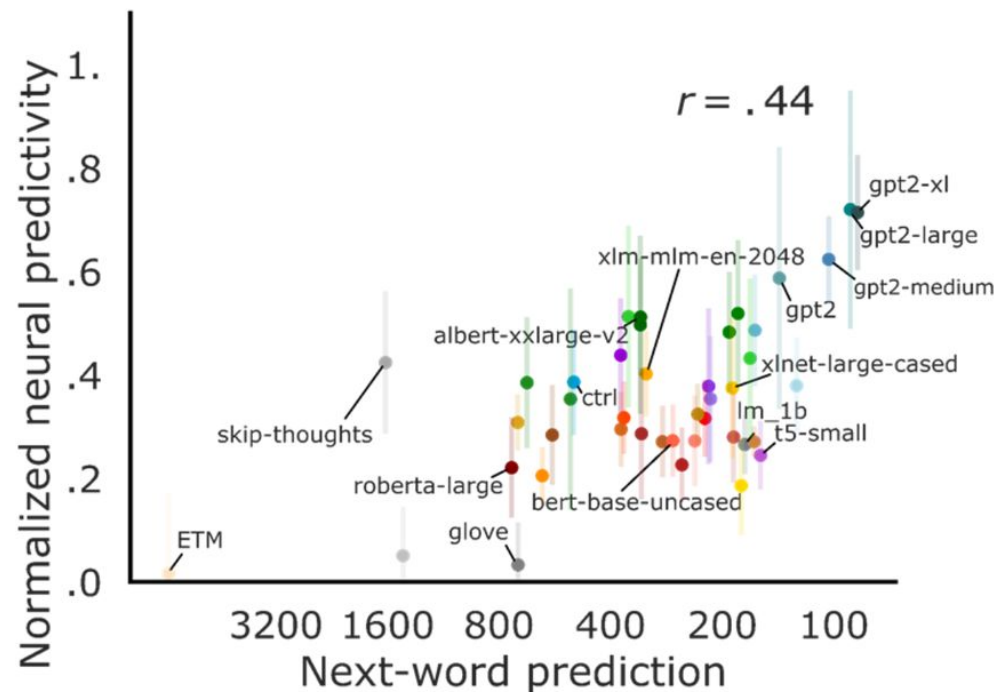
# Syntactic Generalization

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

- 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
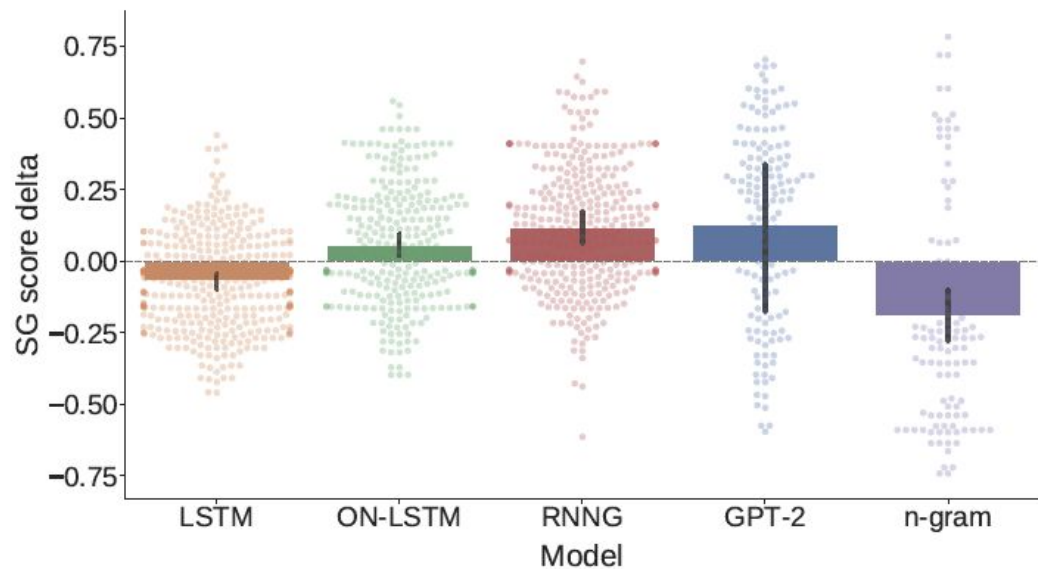
- Compare model class to dataset size

# SG: Perplexity and SG Score

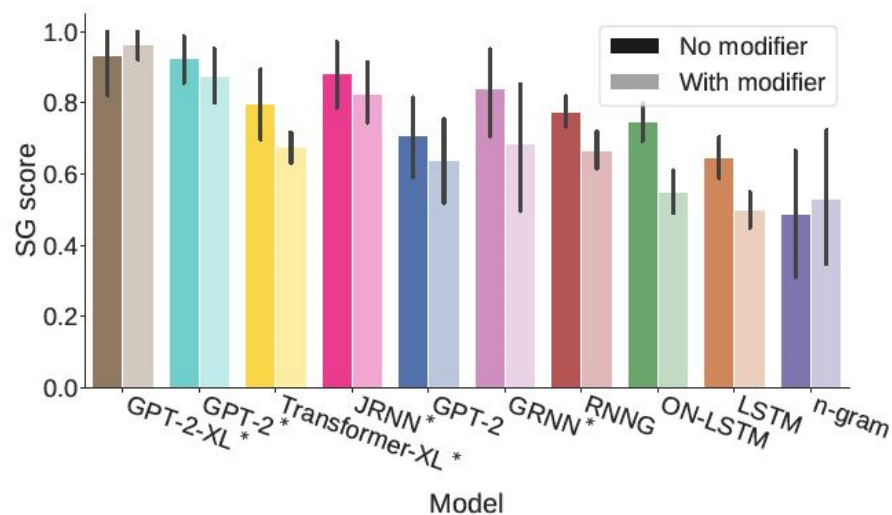# (SG:) Perplexity and Brain-Score

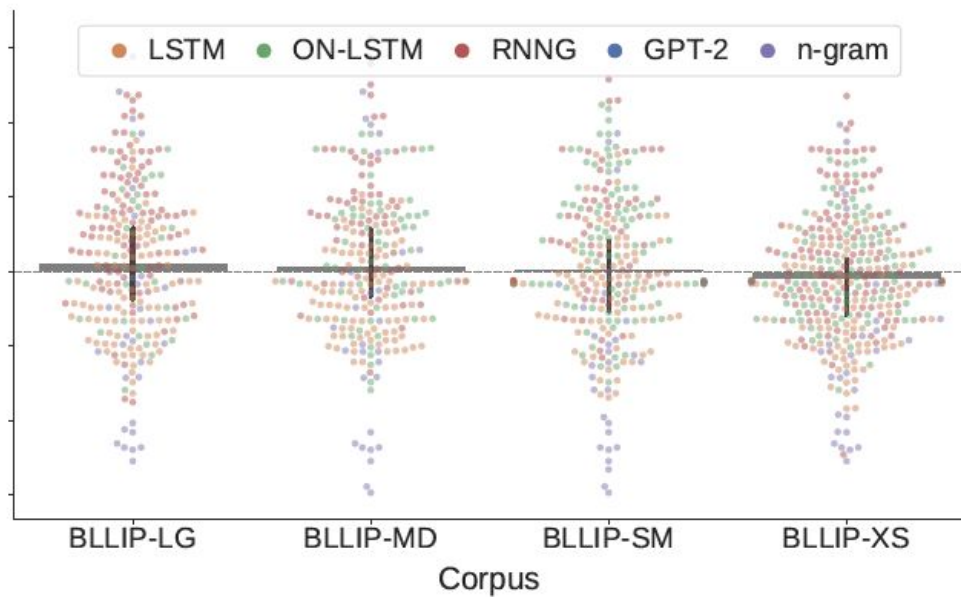# SG: The Influence of Model Architecture

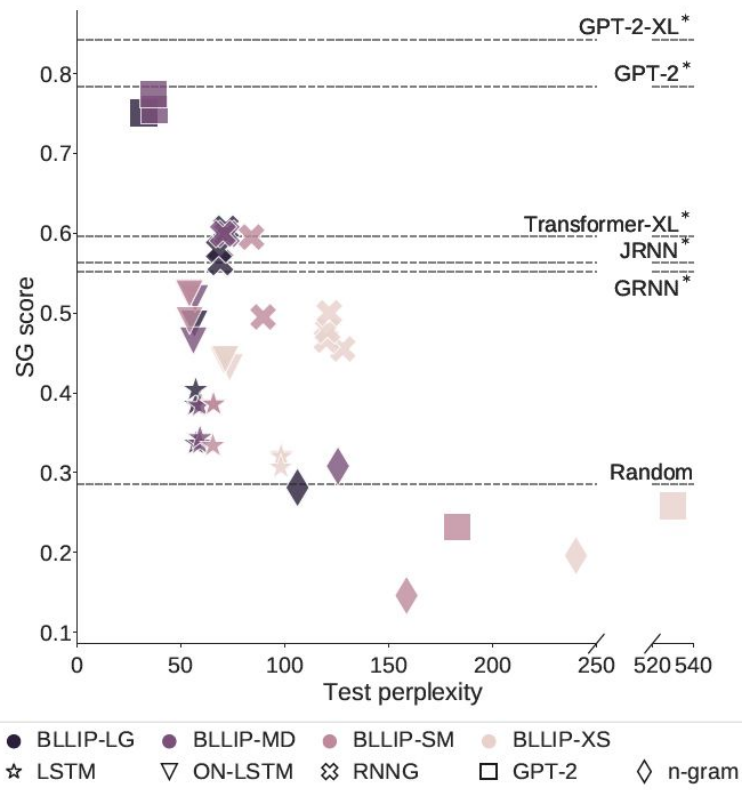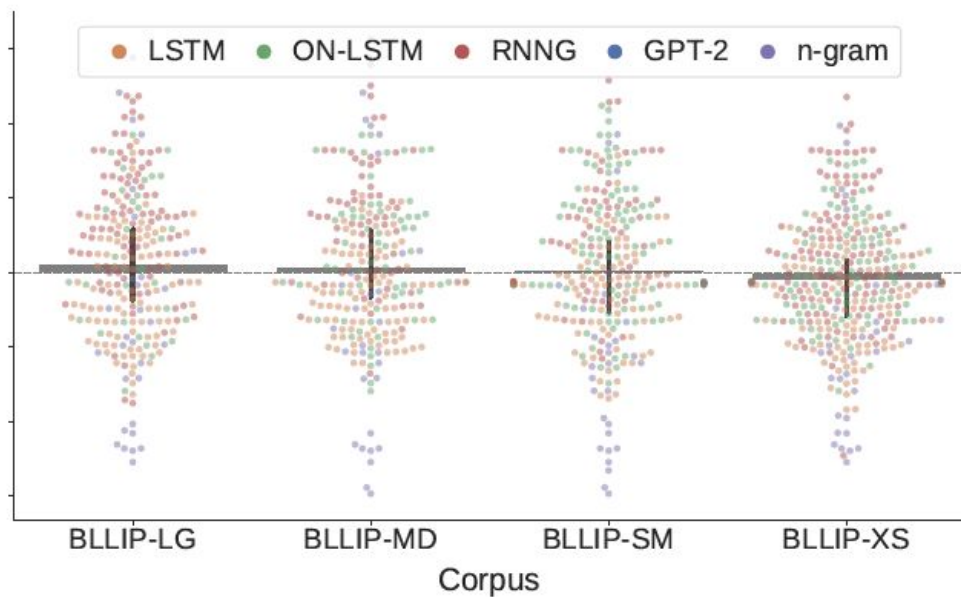# SG: The Influence of Model Architecture

- 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)

# CLUTRR: Motivation and Paradigm

- **C**ompositional **L**anguage **U**nderstanding and **T**ext-based **R**elational **R**easoning
- 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.

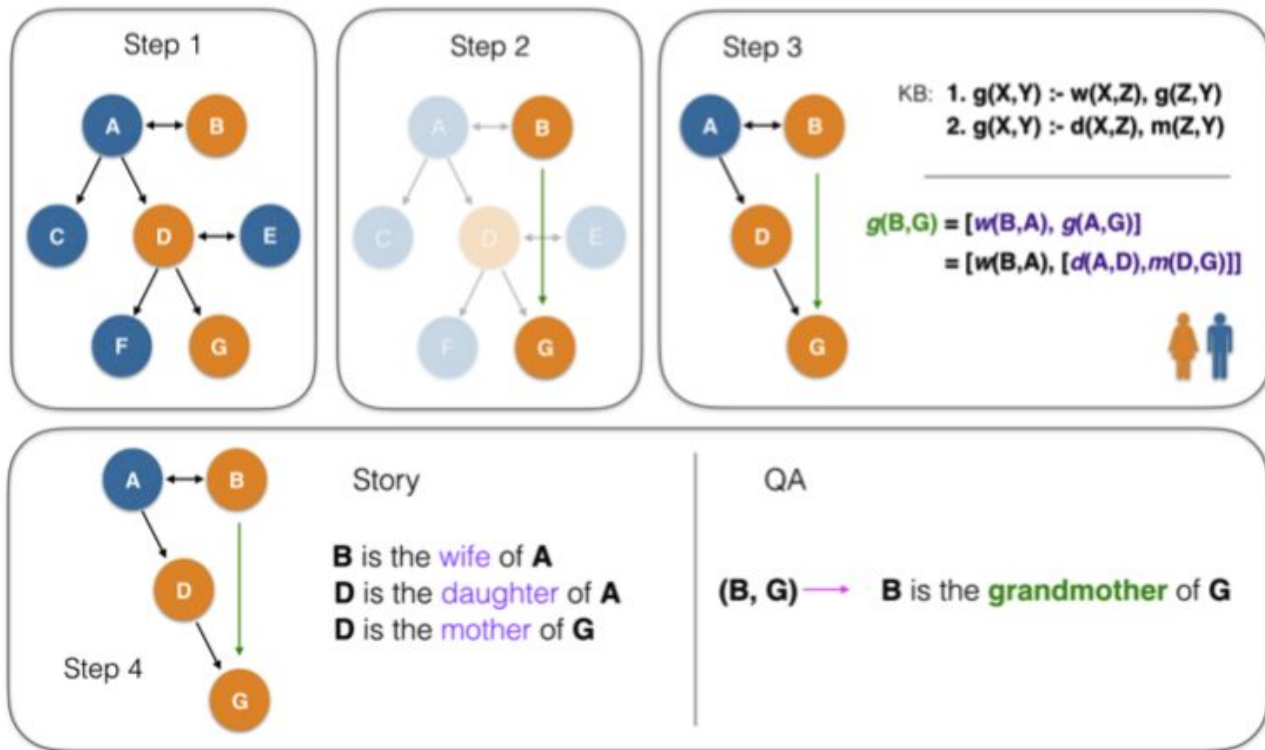Q: How is **Carol** related to **Justin** ?

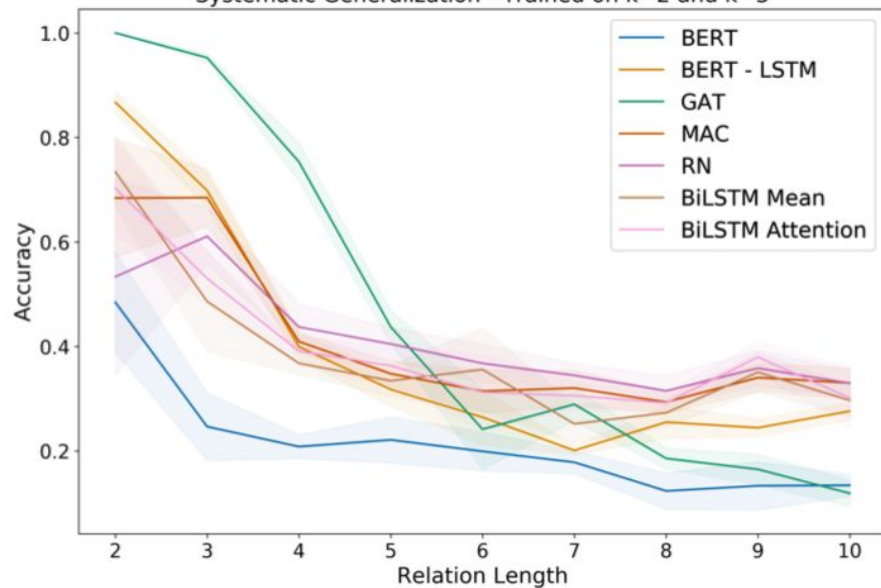A: Carol is the **grandmother** of Justin

# CLUTRR: Motivation and Paradigm

- Productivity
  - mother(mother(mother(Justin))) ~ great grandmother of Justin
- Systematicity
  - Only certain sets allowed with symmetries: son(Justin, Kristin) ~ mother(Kristin, Justin)
- Compositionality
  - son(Justin, Kristin) consists of components

- Memory (compression)
- Children are not exposed to systematic dataset

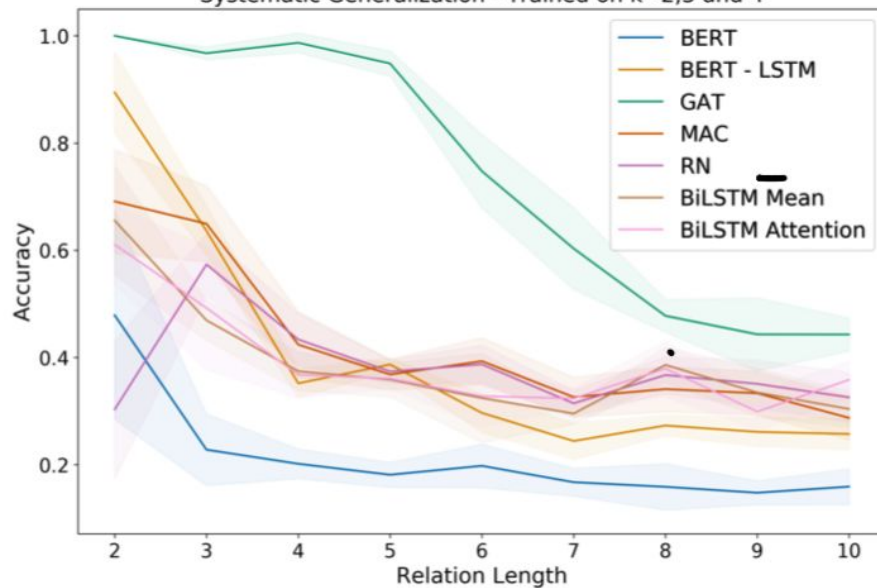# CLUTRR: Dataset Generation & Paradigm

# CLUTRR: Experiment Results

# CLUTRR: Experiment Results

| Models | | Unstructured models (no graph) | | | | | | Structured model (with graph) |
|---|---|---|---|---|---|---|---|---|
| Training | Testing | BiLSTM - Attention | BiLSTM - Mean | RN | MAC | BERT | BERT-LSTM | GAT |
| Clean | Clean | $0.58_{\pm0.05}$ | $0.53_{\pm0.05}$ | $0.49_{\pm0.06}$ | $0.63_{\pm0.08}$ | $0.37_{\pm0.06}$ | $0.67_{\pm0.03}$ | $\mathbf{1.0}_{\pm0.0}$ |
| | Supporting | $\mathbf{0.76}_{\pm0.02}$ | $0.64_{\pm0.22}$ | $0.58_{\pm0.06}$ | $0.71_{\pm0.07}$ | $0.28_{\pm0.1}$ | $0.66_{\pm0.06}$ | $0.24_{\pm0.2}$ |
| | Irrelevant | $0.7_{\pm0.15}$ | $\mathbf{0.76}_{\pm0.02}$ | $0.59_{\pm0.06}$ | $0.69_{\pm0.05}$ | $0.24_{\pm0.08}$ | $0.55_{\pm0.03}$ | $0.51_{\pm0.15}$ |
| | Disconnected | $0.49_{\pm0.05}$ | $0.45_{\pm0.05}$ | $0.5_{\pm0.06}$ | $0.59_{\pm0.05}$ | $0.24_{\pm0.08}$ | $0.5_{\pm0.06}$ | $\mathbf{0.8}_{\pm0.17}$ |
| Supporting | Supporting | $0.67_{\pm0.06}$ | $0.66_{\pm0.07}$ | $0.68_{\pm0.05}$ | $0.65_{\pm0.04}$ | $0.32_{\pm0.09}$ | $0.57_{\pm0.04}$ | $\mathbf{0.98}_{\pm0.01}$ |
| Irrelevant | Irrelevant | $0.51_{\pm0.06}$ | $0.52_{\pm0.06}$ | $0.5_{\pm0.04}$ | $0.56_{\pm0.04}$ | $0.25_{\pm0.06}$ | $0.53_{\pm0.06}$ | $\mathbf{0.93}_{\pm0.01}$ |
| Disconnected | Disconnected | $0.57_{\pm0.07}$ | $0.57_{\pm0.06}$ | $0.45_{\pm0.11}$ | $0.4_{\pm0.1}$ | $0.17_{\pm0.05}$ | $0.47_{\pm0.06}$ | $\mathbf{0.96}_{\pm0.01}$ |
| Average | | $\mathbf{0.61}_{\pm0.08}$ | $0.59_{\pm0.08}$ | $0.54_{\pm0.07}$ | $\mathbf{0.61}_{\pm0.06}$ | $0.30_{\pm0.07}$ | $0.56_{\pm0.05}$ | $\mathbf{0.77}_{\pm0.09}$ |

# CLUTRR: Model Robustness



Supporting Fact | Irrelevant Fact | Disconnected Fact

# 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 | $0.58_{\pm0.05}$ | $0.53_{\pm0.05}$ | $0.49_{\pm0.06}$ | $0.63_{\pm0.08}$ | $0.37_{\pm0.06}$ | $0.67_{\pm0.03}$ | $\mathbf{1.0}_{\pm0.0}$ |
| | Supporting | $\mathbf{0.76}_{\pm0.02}$ | $0.64_{\pm0.22}$ | $0.58_{\pm0.06}$ | $0.71_{\pm0.07}$ | $0.28_{\pm0.1}$ | $0.66_{\pm0.06}$ | $0.24_{\pm0.2}$ |
| | Irrelevant | $0.7_{\pm0.15}$ | $\mathbf{0.76}_{\pm0.02}$ | $0.59_{\pm0.06}$ | $0.69_{\pm0.05}$ | $0.24_{\pm0.08}$ | $0.55_{\pm0.03}$ | $0.51_{\pm0.15}$ |
| | Disconnected | $0.49_{\pm0.05}$ | $0.45_{\pm0.05}$ | $0.5_{\pm0.06}$ | $0.59_{\pm0.05}$ | $0.24_{\pm0.08}$ | $0.5_{\pm0.06}$ | $\mathbf{0.8}_{\pm0.17}$ |
| Supporting | Supporting | $0.67_{\pm0.06}$ | $0.66_{\pm0.07}$ | $0.68_{\pm0.05}$ | $0.65_{\pm0.04}$ | $0.32_{\pm0.09}$ | $0.57_{\pm0.04}$ | $\mathbf{0.98}_{\pm0.01}$ |
| Irrelevant | Irrelevant | $0.51_{\pm0.06}$ | $0.52_{\pm0.06}$ | $0.5_{\pm0.04}$ | $0.56_{\pm0.04}$ | $0.25_{\pm0.06}$ | $0.53_{\pm0.06}$ | $\mathbf{0.93}_{\pm0.01}$ |
| Disconnected | Disconnected | $0.57_{\pm0.07}$ | $0.57_{\pm0.06}$ | $0.45_{\pm0.11}$ | $0.4_{\pm0.1}$ | $0.17_{\pm0.05}$ | $0.47_{\pm0.06}$ | $\mathbf{0.96}_{\pm0.01}$ |
| Average | | $\mathbf{0.61}_{\pm0.08}$ | $0.59_{\pm0.08}$ | $0.54_{\pm0.07}$ | $\mathbf{0.61}_{\pm0.06}$ | $0.30_{\pm0.07}$ | $0.56_{\pm0.05}$ | $\mathbf{0.77}_{\pm0.09}$ |

# Future work & Perspectives

- Sub-word tokenization
- Common-sense reasoning
- Abstractions as probabilistic
-

# References

Brysbaert, M., Stevens, M., Mandera, 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.