

# Qualitative Probing of Deep Contextual Models: We need your help!

**Kyle Richardson**

Allen Institute for Artificial Intelligence (AI2)

December 2020

# Probing Natural Language Understanding (NLU) Models

- ▶ **Probing**: understanding the strengths/weaknesses of models ; measuring model competence qualitatively; **behavioral (input/output) testing**.

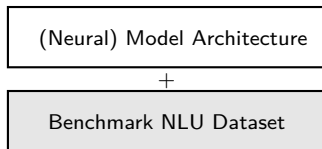


# Building NLU Models: Standard Picture

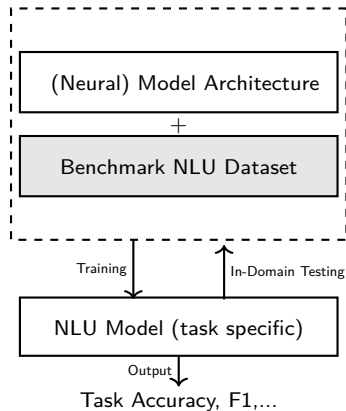
# Building NLU Models: Standard Picture

(Neural) Model Architecture

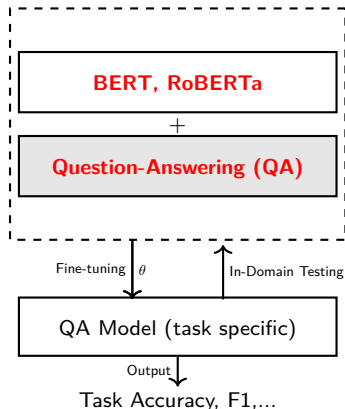
# Building NLU Models: Standard Picture



# Building NLU Models: Standard Picture

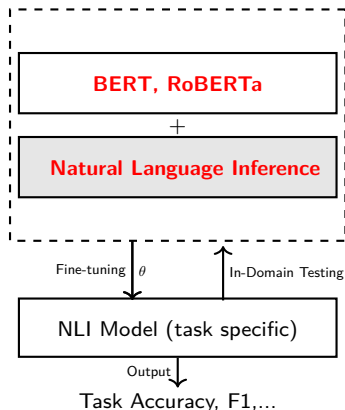


# Building NLU Models: Standard Picture



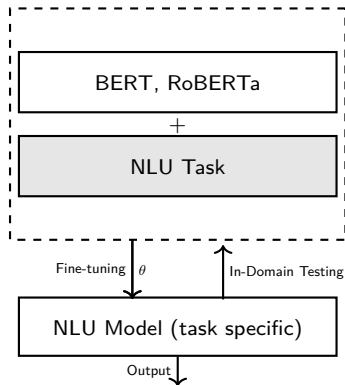
Multiple-Choice QA (ARC Benchmark)	
<b>Question</b>	<i>Which property of a mineral can be determined just by looking at it?</i>
<b>Answers</b>	(A) <u>luster</u> (B) <u>mass</u> (C) <u>weight</u> (D) <u>hardness</u>
	correct answer      distractor 1      distractor 2      distractor 3

# Building NLU Models: Standard Picture



Natural Language Inference (SNLI benchmark)	
<b>Sen1</b>	<i>A soccer game with multiple males playing.</i>
<b>Sen2</b>	<i>Some men are playing a sport.</i>
<b>Label</b>	Yes/Entailment

# Qualitative Analysis of Models



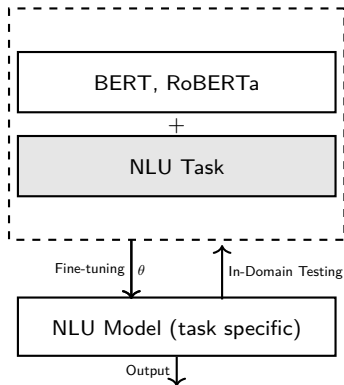
## Desiderata

Does my model know about

Taxonomic relations, definitions, synonymy,  
robust to perturbations/consistent, ....?

- **Why?** Models sometimes do the right things for the wrong reasons ; exploit biases (Gururangan et al., 2018); **model/bug repair**.

# Qualitative Analysis of Models



## Desiderata

Does my model know about

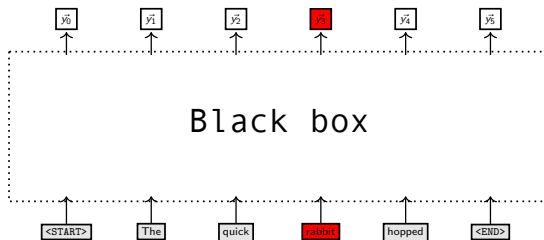
Taxonomic relations, definitions, synonymy,  
robust to perturbations/consistent, ....?

- **Bigger Issue** (not often discussed) Unclear how linguists, logicians, people working on classical AI fit into this picture; facilitate collab.



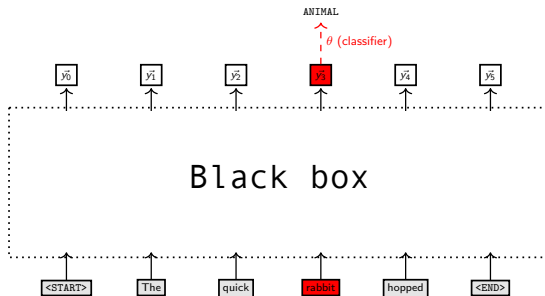
# Contextual Models: A Cursory Overview

- **Role:** assign continuous (non-symbolic) vector representations  $y \in \mathbb{R}$  to inputs based on their meaning in each instance; deep neural networks.



# Contextual Models: A Cursory Overview

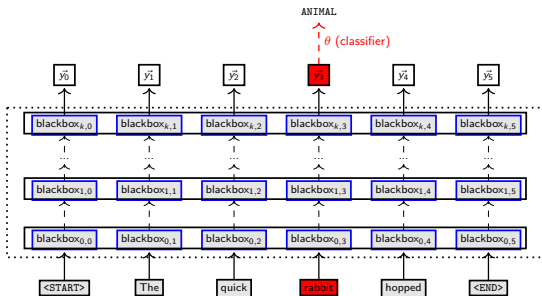
- **Role:** assign continuous (non-symbolic) vector representations  $y \in \mathbb{R}$  to inputs based on their meaning in each instance; deep neural networks.



- **Words as Vectors:** Give models considerable power; word/concept similarity reduces to vector similarity, e.g.,  $\text{SIMILARITY}(\overrightarrow{\text{rabbit}}, \overrightarrow{\text{bunny}})$ .

# Contextual Models: A Cursory Overview

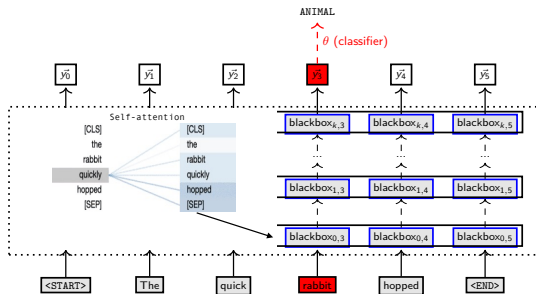
- **Role:** assign continuous (non-symbolic) vector representations  $y \in \mathbb{R}$  to inputs based on their meaning in each instance; deep neural networks.



- **Development 1:** New architecture, **Transformers** (Vaswani et al., 2017), dispense with recurrent sequential structures, self-attention.

# Contextual Models: A Cursory Overview

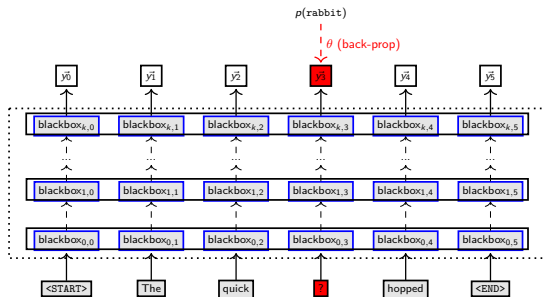
- **Role:** assign continuous (non-symbolic) vector representations  $y \in \mathbb{R}$  to inputs based on their meaning in each instance; deep neural networks.



- **Development 1:** New architecture, **Transformers** (Vaswani et al., 2017), dispense with recurrent sequential structures, self-attention.

# Contextual Models: A Cursory Overview

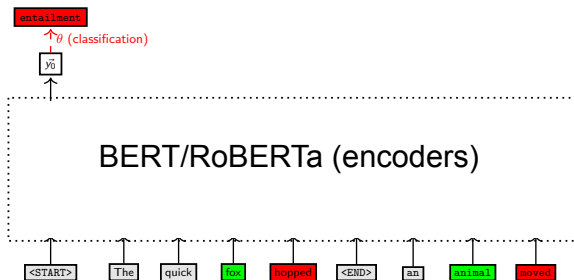
- **Role:** assign continuous (non-symbolic) vector representations  $y \in \mathbb{R}$  to inputs based on their meaning in each instance; deep neural networks.



- **Development 2: Model pre-training:** Have the model read the internet (terabytes of data) and learn by solving word completion (cloze) tasks.

# Contextual Models: A Cursory Overview

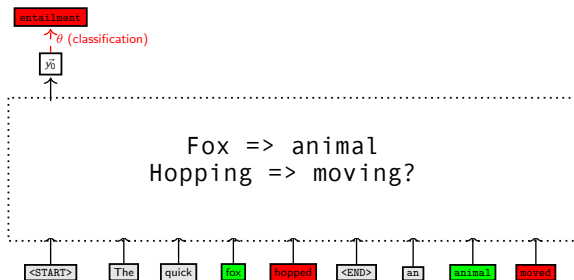
- **Role:** assign continuous (non-symbolic) vector representations  $y \in \mathbb{R}$  to inputs based on their meaning in each instance; deep neural networks.



- **Fine-tuning:** Training models on smaller customized tasks, exploiting pre-trained knowledge. Pioneered in [Devlin et al. \(2018\)](#).

# Contextual Models: A Cursory Overview

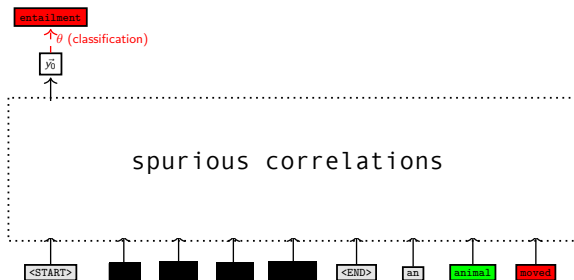
- **Role:** assign continuous (non-symbolic) vector representations  $y \in \mathbb{R}$  to inputs based on their meaning in each instance; deep neural networks.



- Are these models actually knowledgeable, or just exploiting tricks and systematic biases in data? [Gururangan et al. \(2018\)](#)

# Contextual Models: A Cursory Overview

- **Role:** assign continuous (non-symbolic) vector representations  $y \in \mathbb{R}$  to inputs based on their meaning in each instance; deep neural networks.



- Are these models actually knowledgeable, or just exploiting tricks and systematic biases in data? [Gururangan et al. \(2018\)](#)



# Requirements for Demonstrating knowledge

- ▶ QA in the **Science domain**, well studied qualitatively ([Clark et al., 2018](#); [Boratko et al., 2018](#)), though anecdotal and post-hoc .

# Requirements for Demonstrating knowledge

- ▶ QA in the **Science domain**, well studied qualitatively ([Clark et al., 2018](#); [Boratko et al., 2018](#)), though anecdotal and post-hoc .

ARC Challenge ( <a href="#">Clark et al., 2018</a> )	
<b>Question</b>	What is a worldwide increase in temperature <u>called</u> ? <i>Definition</i>
<b>Answers</b>	(A) greenhouse effect (B) global warming (C) ozone depletion (D) solar heating.
<b>Knowledge:</b>	DEF( <i>global warming, worldwide increase in...</i> )

# Requirements for Demonstrating knowledge

- QA in the **Science domain**, well studied qualitatively ([Clark et al., 2018](#); [Boratko et al., 2018](#)), though anecdotal and post-hoc .

ARC Challenge ( <a href="#">Clark et al., 2018</a> )	
Question	What is a worldwide increase in temperature <u>called</u> ? <i>Definition</i>
Answers	(A) greenhouse effect (B) <b>global warming</b> (C) ozone depletion (D) solar heating.
Knowledge:	DEF( <b>global warming, worldwide increase in...</b> )

OpenBookQA ( <a href="#">Mihaylov et al., 2018</a> )	
Question	Which of the following <u>is a type of learned behavior</u> ? <i>ISA reasoning</i>
Answers	(A) <b>cooking</b> (B) thinking (C) hearing (D) breathing
Knowledge:	ISA( <b>cooking, learned behavior</b> )

# Requirements for Demonstrating knowledge

- QA in the **Science domain**, well studied qualitatively (Clark et al., 2018; Boratko et al., 2018), though anecdotal and post-hoc .

ARC Challenge (Clark et al., 2018)	
Question	What is a worldwide increase in temperature <u>called</u> ? <i>Definition</i>
Answers	(A) greenhouse effect (B) global warming (C) ozone depletion (D) solar heating.
Knowledge:	DEF(global warming,worldwide increase in...)

OpenBookQA (Mihaylov et al., 2018)	
Question	Which of the following <u>is a type of learned behavior</u> ? <i>ISA reasoning</i>
Answers	(A) cooking (B) thinking (C) hearing (D) breathing
Knowledge:	ISA(cooking,learned behavior)

*Do models truly possess the basic knowledge/reasoning skills we think they do?* Hard to say without **specialized tests**.

# Requirements for Demonstrating Knowledge

- QA in the **Science domain**, well studied qualitatively (Clark et al., 2018; Boratko et al., 2018), though anecdotal and post-hoc .

ARC Challenge (Clark et al., 2018)	
Question	What is a worldwide increase in temperature <u>called</u> ? <i>Definition</i>
Answers	(A) greenhouse effect (B) global warming (C) ozone depletion (D) solar heating.
Knowledge:	DEF(global warming,worldwide increase in...)

OpenBookQA (Mihaylov et al., 2018)	
Question	Which of the following <u>is a type of</u> learned behavior? <i>ISA reasoning</i>
Answers	(A) cooking (B) thinking (C) hearing (D) breathing
Knowledge:	ISA(cooking,learned behavior)

To demonstrate competence a **model should:**

1. have knowledge across *many concepts*;
2. be robust to *perturbations*
3. *and varying levels of reasoning complexity* .

# Requirements for Demonstrating Knowledge

- QA in the **Science domain**, well studied qualitatively (Clark et al., 2018; Boratko et al., 2018), though anecdotal and post-hoc .

ARC Challenge (Clark et al., 2018)	
Question	What is <i>the thinning of Earth's upper atmosphere</i> <u>called</u> ? <i>Definition</i>
Answers	(A) greenhouse effect (B) global warming (C) <i>ozone depletion</i> (D) solar heating.
Knowledge:	DEF( <i>ozone depletion, thinning of the Earth's ...</i> )

OpenBookQA (Mihaylov et al., 2018)	
Question	Which of the following <u>is a type of</u> learned behavior? <i>ISA reasoning</i>
Answers	(A) <i>cooking</i> (B) thinking (C) hearing (D) breathing
Knowledge:	ISA( <i>cooking, learned behavior</i> )

To demonstrate competence a **model should:**

1. have knowledge across *many concepts*;
2. be robust to *perturbations*
3. *and varying levels of reasoning complexity* .

# Requirements for Demonstrating Knowledge

- QA in the **Science domain**, well studied qualitatively (Clark et al., 2018; Boratko et al., 2018), though anecdotal and post-hoc .

ARC Challenge (Clark et al., 2018)	
Question	What is <i>the thinning of Earth's upper atmosphere</i> <u>called</u> ? <i>Definition</i>
Answers	(A) greenhouse effect (B) global warming (C) <i>ozone depletion</i> (D) solar heating.
Knowledge:	DEF( <i>ozone depletion, thinning of the Earth's ...</i> )

OpenBookQA (Mihaylov et al., 2018)	
Question	Which of the following is a <u>type form</u> of <i>learned</i> behavior? <i>ISA reasoning</i>
Answers	(A) <i>cooking</i> (B) thinking (C) hearing (D) <i>breathing eating</i>
Knowledge:	ISA( <i>cooking, learned behavior</i> )

To demonstrate competence a **model should:**

1. have knowledge across *many concepts*;
2. be robust to *perturbations*
3. *and varying levels of reasoning complexity* .

# Requirements for Demonstrate Knowledge

- ▶ When demonstrating knowledge, want to consider the extreme cases with considerable complexity; can result in pedantic English.

sen1	sen2	Label
Mitchell is as tall as Fred, Fred is as tall as Karl, Karl is as tall as Jon, Jon is as tall as Darryl, Darryl is as tall as Theodore, Theodore is as tall as Calvin, Calvin is as tall as Eddie , Eddie is as tall as Philip , Philip is taller than Travis	Calvin is taller than Travis .	Entailment
A bat with a strong odor did not hit several dogs	A bat with a strong smell did not hit many poodles	Entailment



# Requirements for Demonstrate Knowledge

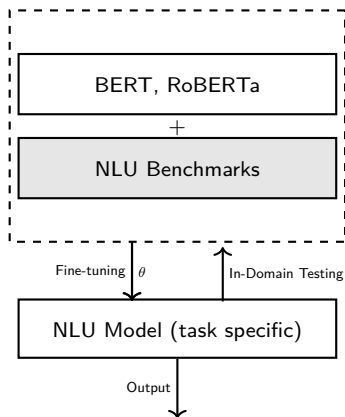
- ▶ When demonstrating knowledge, want to consider the extreme cases with considerable complexity; can result in pedantic English.

sen1	sen2	Label
Mitchell is as tall as Fred, Fred is as tall as Karl, Karl is as tall as Jon, Jon is as tall as Darryl, Darryl is as tall as Theodore, Theodore is as tall as Calvin, Calvin is as tall as Eddie , Eddie is as tall as Philip , Philip is taller than Travis	Calvin is taller than Travis .	Entailment
A bat with a strong odor did not hit several dogs	A bat with a strong smell did not hit many poodles	Entailment

- ▶ **Other robustness measures:** out of domain testing, lexical diversity ([Rozen et al., 2019](#)).

# Diagnostic Tasks for NLU

- ▶ unit testing (Ribeiro et al., 2020), challenge tasks/stress tests (**task specific**) (Naik et al., 2018; Glockner et al., 2018), *inter alia*.

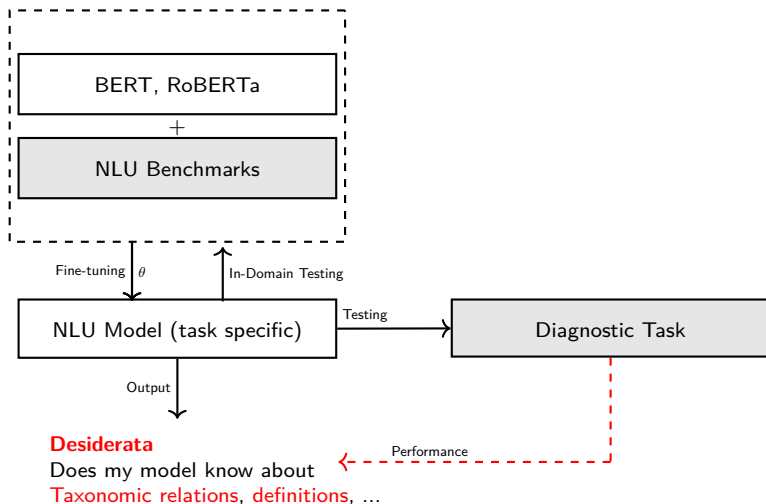


## Desiderata

Does my model know about  
Taxonomic relations, definitions, ...

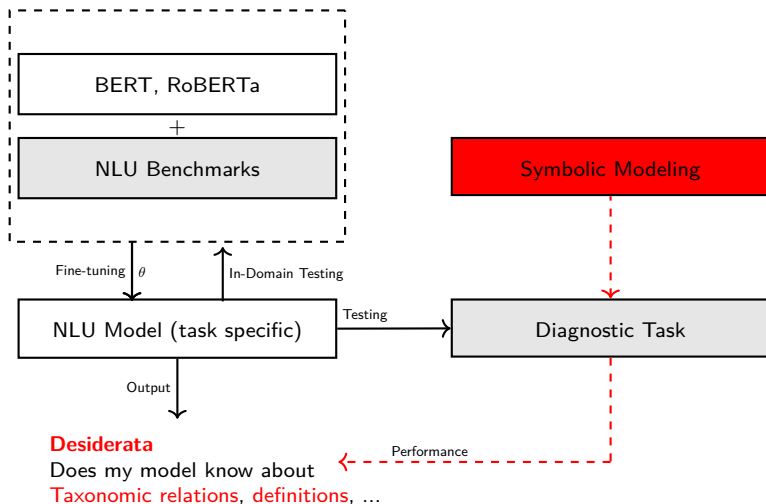
# Diagnostic Tasks for NLU

- ▶ unit testing (Ribeiro et al., 2020), challenge tasks/stress tests (**task specific**) (Naik et al., 2018; Glockner et al., 2018), *inter alia*.



# Diagnostic Tasks for NLU

- ▶ unit testing (Ribeiro et al., 2020), challenge tasks/stress tests (**task specific**) (Naik et al., 2018; Glockner et al., 2018), *inter alia*.



# <Building Diagnostic Tasks>

(3 Example Studies)

# Diagnostic Tasks via Expert Knowledge (Richardson and Sabharwal, 2020)[TACL]

- ▶ A model should 1. have knowledge across many concepts ; 2. robust to perturbations ; 3. varying complexity .

# Diagnostic Tasks via Expert Knowledge (Richardson and Sabharwal, 2020)[TACL]

- ▶ A ~~model~~ **dataset** should 1. ~~have~~ **test** knowledge across many concepts ;  
2. ~~robust to~~ **have** perturbations ; 3. ~~varying~~ **controlled** complexity .

**Assumption:** we can demonstrate that models exhibit these properties by testing them on data that has these properties..

# Diagnostic Tasks via Expert Knowledge (Richardson and Sabharwal, 2020)[TACL]

- ▶ A ~~model~~ **dataset** should 1. ~~have~~ **test** knowledge across many concepts ;  
2. ~~robust to~~ **have** perturbations ; 3. ~~varying~~ **controlled** complexity .

**Assumption:** we can demonstrate that models exhibit these properties by testing them on data that has these properties..

Expert Knowledge (KBs, lexical ontologies)

Arg1	Arg2	REL	EX
nestle.v	position comfort- ably	DEF	The baby nestled her head..
elude.v	escape.v	ISA	The thief eluded po- lice...
trouser.n	consumer good.n	ISA	The man bought trousers..
poet.n	writer.n	ISA	...
...	...	...	...



# Diagnostic Tasks via Expert Knowledge (Richardson and Sabharwal, 2020)[TACL]

- ▶ A ~~model~~ **dataset** should 1. ~~have~~ **test** knowledge across many concepts ;
- 2. ~~robust to~~ **have** perturbations ; 3. ~~varying~~ **controlled** complexity .

**Assumption:** we can demonstrate that models exhibit these properties by testing them on data that has these properties..

Expert Knowledge (KBs, lexical ontologies)

Arg1	Arg2	REL	EX
nestle.v	position comfort- ably	DEF	The baby nestled her head..
elude.v	escape.v	ISA	The thief eluded po- lice...
trouser.n	consumer good.n	ISA	The man bought trousers..
poet.n	writer.n	ISA	...
...	...	...	...

templates

Probing Questions

Question	Answer	Test
Given 'The baby nes- tled her head', <b>nes- tled</b> is defined as	position comfort- ably	def
In 'we had to spell our name for the police', <b>spell</b> is a type of	recite event	isa
In the context, 'the poet published his new poem', <b>poet</b> is best defined as	a writer of poems....	def

# Diagnostic Tasks via Expert Knowledge (Richardson and Sabharwal, 2020)[TACL]

- ▶ A ~~model~~ **dataset** should 1. ~~have~~ **test** knowledge across many concepts ;
- 2. ~~robust to~~ **have** perturbations ; 3. ~~varying~~ **controlled** complexity .

**Assumption:** we can demonstrate that models exhibit these properties by testing them on data that has these properties..

Expert Knowledge (KBs, lexical ontologies)

Arg1	Arg2	REL	EX
nestle.v	position comfort- ably	DEF	The baby nestled her head..
elude.v	escape.v	ISA	The thief eluded po- lice...
trouser.n	consumer good.n	ISA	The man bought trousers..
poet.n	writer.n	ISA	...
...	...	...	...

templates

Probing Questions

Question	Answer	Test
Given 'The baby nes- tled her head', nes- tled is defined as	position comfort- ably	def
In 'we had to spell our name for the police', spell is a type of	recite event	isa
In the context, 'the poet published his new poem', poet is best defined as	a writer of poems....	def

distractor assignment/taxonomic constraints

Diagnostic Task

# Diagnostic Tasks via Expert Knowledge (Richardson and Sabharwal, 2020)[TACL]

- ▶ A ~~model~~ **dataset** should 1. ~~have~~ **test** knowledge across many concepts ;  
2. ~~robust to~~ **have** perturbations ; 3. ~~varying~~ **controlled** complexity .

**Assumption:** we can demonstrate that models exhibit these properties by testing them on data that has these properties..

Expert Knowledge (KBs, lexical ontologies)

Arg1	Arg2	REL	EX
nestle.v	position comfort- ably	DEF	The baby nestled her head..
elude.v	escape.v	ISA	The thief eluded po- lice...
trouser.n	consumer good.n	ISA	The man bought trousers..
poet.n	writer.n	ISA	...
...	...	...	...

templates

Probing Questions

Question	Answer	Test
Given 'The baby nes- tled her head', nes- tled is defined as	position comfort- ably	def
In 'we had to spell our name for the police', spell is a type of	recite event	isa
In the context, 'the poet published his new poem', poet is best defined as	a writer of poems....	def

distractor assignment/taxonomic constraints

**Diagnostic Task**

**Meta-level QA:** Asking questions about abstract knowledge; many concepts (1. ✓); controlled templates/distractor complexity (2. ✓ 3. ✓)

# Diagnostic Tasks via Expert Knowledge (Richardson and Sabharwal, 2020)[TACL]

- ▶ A ~~model~~ **dataset** should 1. ~~have~~ **test** knowledge across many concepts ;  
2. ~~robust to~~ **have** perturbations ; 3. ~~varying~~ **controlled** complexity .

**Assumption:** we can demonstrate that models exhibit these properties by testing them on data that has these properties..

Expert Knowledge (KBs, lexical ontologies)

Arg1	Arg2	REL	EX
nestle.v	position comfort- ably	DEF	The baby nestled her head..
elude.v	escape.v	ISA	The thief eluded po- lice...
trouser.n	consumer good.n	ISA	The man bought trousers..
poet.n	writer.n	ISA	...
...	...	...	...

templates

Probing Questions

Question	Answer	Test
Given 'The baby nes- tled her head', <b>nes- tled</b> is defined as	position comfort- ably	def
In 'we had to spell our name for the police', <b>spell</b> is a type of	recite event	isa
In the context, 'the poet published his new poem', <b>poet</b> is best defined as	a writer of poems....	def

distractor assignment/taxonomic constraints

**Diagnostic Task**

**Trade-offs:** KBs tends to be noisy; dealt by synthesizing large amount of data, contextualizing questions, gold test annotation (where needed).

# Example QA Diagnostics

- **Resources:** WordNet, GCIDE dictionary; **5 individual tasks:** Definitions, Synonymy, Hypernymy (ISA), and Hyponymy (ISA), WordSense.
- WordNet tasks involve  $\sim 30k$  atomic concepts, exhaustive combinations of distractors.

Probe	Example
Definitions + Word Sense	In the sentence <i>The baby nestled her head</i> , the word <i>nestled</i> is best defined as (A) <u>position comfortably</u> (B) <u>put in a certain place</u> (C) <u>a type of fish</u> ... <i>correct answer</i> <i>hard/close distractor</i> <i>easy/random distractor</i>
Hypernymy (ISA)	In <i>The thief eluded the police</i> , the word of concept <i>eluded</i> is best described as (A) ... (B) <u>an escape event, defined as ...</u> (C) ... <i>correct answer</i>
Hyponymy (ISA)	Given the context <i>They awaited her arrival</i> , which of the following is a specific type of <i>arrival</i> (A) <u>driving a car</u> (B) <u>crash landing, defined as .....</u> related concept <i>correct answer</i>
Synonymy	Which set of words best corresponds to the definition of <i>a grammatical category in inflected languages...</i> (A) <u>gender</u> (B) ... <i>correct answer</i>

# Semantic Fragments for NLI

Richardson et al. (2020)[AAAI]

- ▶ **Semantic Fragment:** subset of language *equipped with semantics* which *translate into some formal system* ... ([Pratt-Hartmann, 2004](#))

# Semantic Fragments for NLI Richardson et al. (2020)[AAAI]

- **Semantic Fragment:** subset of language *equipped with semantics* which *translate into some formal system* ... ([Pratt-Hartmann, 2004](#))

**Formal Specification of Facts about Quantifiers** ([van Benthem \(1986\)](#))

<u>all</u> $X \ Y$	$\models$	<u>all</u> $X' \ Y'$ , s.t. $X' \leq X, Y \leq Y'$
<u>some</u> $X \ Y$	$\models$	<u>some</u> $X' \ Y'$ , s.t. $X \leq X', \dots$
<u>exactly</u> $N \ X \ ..$	$\models$	....

# Semantic Fragments for NLI Richardson et al. (2020)[AAAI]

- **Semantic Fragment:** subset of language *equipped with semantics* which *translate into some formal system* ... (Pratt-Hartmann, 2004)

**Formal Specification of Facts about Quantifiers** (van Benthem (1986))

<u>all</u> $X \ Y$	$\models$	<u>all</u> $X' \ Y'$ , s.t. $X' \leq X, Y \leq Y'$
<u>some</u> $X \ Y$	$\models$	<u>some</u> $X' \ Y'$ , s.t. $X \leq X', \dots$
<u>exactly</u> $N \ X \ ..$	$\models$	....

**Example Semantic Fragment**

↓ symbolic model+generator+lexicon

*All dogs ran*  $\models$  *All small dogs ran*, *All furry dogs barked*  $\models$  *All animals barked*, *Some dog ran*  $\models$  *Some animal moved*, ...  $\uparrow$



# Semantic Fragments for NLI Richardson et al. (2020)[AAAI]

- **Semantic Fragment:** subset of language *equipped with semantics* which *translate into some formal system* ... (Pratt-Hartmann, 2004)

**Formal Specification of Facts about Quantifiers** (van Benthem (1986))

<u>all</u> X Y	$\models$	<u>all</u> X' Y', s.t. $X' \leq X, Y \leq Y'$
<u>some</u> X Y	$\models$	<u>some</u> X' Y', s.t. $X \leq X', \dots$
<u>exactly</u> N X ..	$\models$	....

**Example Semantic Fragment**

↓ symbolic model+generator+lexicon

*All dogs ran*  $\models$  *All small dogs ran*, *All furry dogs barked*  $\models$  *All animals barked*, *Some dog ran*  $\models$  *Some animal moved*, ...  $\uparrow$

↓ NLI format, standard splitting, ...

Diagnostic Task

# Semantic Fragments for NLI Richardson et al. (2020)[AAAI]

- **Semantic Fragment:** subset of language *equipped with semantics* which *translate into some formal system* ... (Pratt-Hartmann, 2004)

**Formal Specification of Facts about Quantifiers** (van Benthem (1986))

<u>all</u> X Y	$\models$	<u>all</u> X' Y', s.t. $X' \leq X, Y \leq Y'$
<u>some</u> X Y	$\models$	<u>some</u> X' Y', s.t. $X \leq X', \dots$
<u>exactly</u> N X ..	$\models$	....

**Example Semantic Fragment**

↓ symbolic model+generator+lexicon

All dogs ran  $\models$  All small dogs ran, All furry dogs barked  $\models$  All animals barked, Some dog ran  $\models$  Some animal moved,...  $\uparrow$

↓ NLI format, standard splitting,...

Diagnostic Task

**Non-standard in NLP:** Using symbolic models (vs. humans) to elicit data; standard tool in linguistics (Montague (1973)).

# Semantic Fragments for NLI Richardson et al. (2020)[AAAI]

- **7 Tasks:** elementary logic (e.g., boolean algebra, quantification, conditionals) and monotonicity reasoning

Fragments	Example (premise,label,hypothesis)
Negation	<i>Laurie has only visited Nephi, Marion has only visited Calistoga.</i> CONTRADICTION <i>Laurie didn't visit Nephi</i>
Boolean	<i>Travis, Arthur, Henry and Dan have only visited Georgia</i> ENTAILMENT <i>Dan didn't visit Rwanda</i>
Quantifier	<i>Everyone has visited every place</i> NEUTRAL <i>Virgil didn't visit Barry</i>
Counting	<i>Nellie has visited Carrie, Billie, John, Mike, Thomas, Mark, ..., and Arthur.</i> ENTAILMENT <i>Nellie has visited more than 10 people.</i>
Conditionals	<i>Francisco has visited Potsdam and if Francisco has visited Potsdam then Tyrone has visited Pampa</i> ENTAILMENT <i>Tyrone has visited Pampa.</i>
Comparatives	<i>John is taller than Gordon and Erik..., and Mitchell is as tall as John</i> NEUTRAL <i>Erik is taller than Gordon.</i>
Monotonicity	<i>All black mammals saw exactly 5 stallions who danced</i> ENTAILMENT <i>A brown or black poodle saw exactly 5 stallions who danced</i>

# Semantic Fragments for NLI Richardson et al. (2020)[AAAI]

- **7 Tasks:** elementary logic (e.g., boolean algebra, quantification, conditionals) and monotonicity reasoning

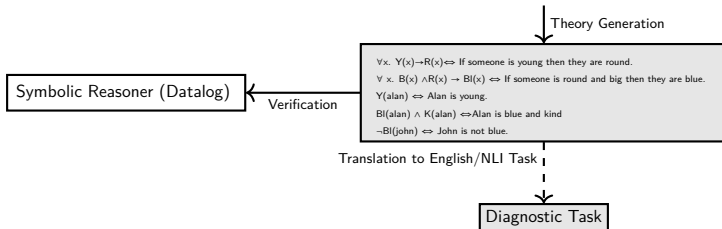
Fragments	Example (premise,label,hypothesis)
Negation	<i>Laurie has only visited Nephi, Marion has only visited Calistoga.</i> CONTRADICTION <i>Laurie didn't visit Nephi</i>
Boolean	<i>Travis, Arthur, Henry and Dan have only visited Georgia</i> ENTAILMENT <i>Dan didn't visit Rwanda</i>
Quantifier	<i>Everyone has visited every place</i> NEUTRAL <i>Virgil didn't visit Barry</i>
Counting	<i>Nellie has visited Carrie, Billie, John, Mike, Thomas, Mark, ..., and Arthur.</i> ENTAILMENT <i>Nellie has visited more than 10 people.</i>
Conditionals	<i>Francisco has visited Potsdam and if Francisco has visited Potsdam then Tyrone has visited Pampa</i> ENTAILMENT <i>Tyrone has visited Pampa.</i>
Comparatives	<i>John is taller than Gordon and Erik..., and Mitchell is as tall as John</i> NEUTRAL <i>Erik is taller than Gordon.</i>
Monotonicity	<i>All black mammals saw exactly 5 stallions who danced</i> ENTAILMENT <i>A brown or black poodle saw exactly 5 stallions who danced</i>

Done in collaboration with logicians and linguists (Indiana University);  
generated using simple templates , formal grammars.

# Rule Taker: Training Models to do Formalized Reasoning.

(Clark et al., 2020)[IJCAI]

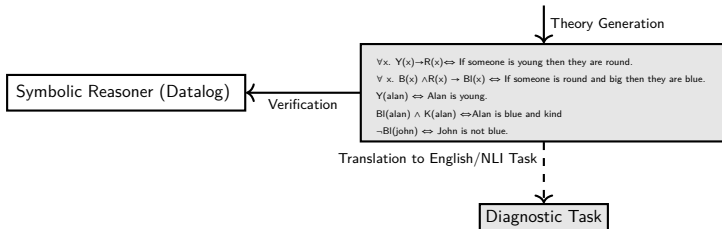
- **Idea:** Synthesize deductively valid entailment data (theories and queries) with the help of symbolic theorem prover; render as English.



# Rule Taker: Training Models to do Formalized Reasoning.

(Clark et al., 2020)[IJCAI]

- **Idea:** Synthesize **deductively valid entailment data** (theories and queries) with the help of symbolic theorem prover; render as English.



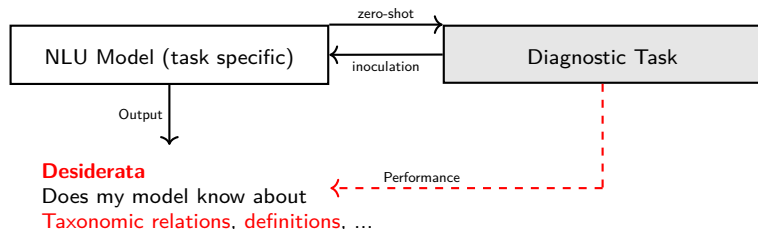
**Bigger idea:** demonstrating model correctness can be achieved by testing on data that is *correct by construction*.

- **Components:** reasoning depth, vocabulary overlap/mismatch.

# Probing Methodology and Experiments (QA+NLI Tasks)

- ▶ Trained single models on standard benchmarks; **Ask the following empirical questions:**

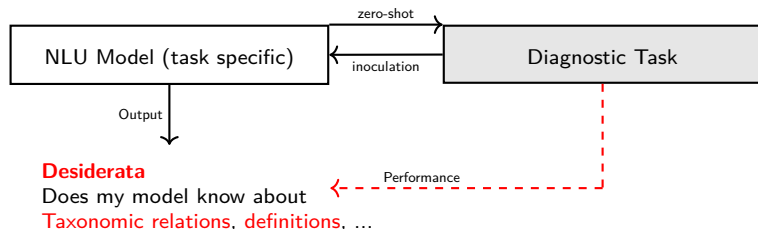
1. How well do benchmark models perform on each *individual* probing on diagnostic task without specialized training (**zero-shot**)?
2. How well models perform after a small amount of additional training on probes (**inoculation** (Liu et al., 2019))?



# Probing Methodology and Experiments (QA+NLI Tasks)

- ▶ Trained single models on standard benchmarks; **Ask the following empirical questions:**

1. How well do benchmark models perform on each *individual* probing on diagnostic task without specialized training (**zero-shot**)?
2. How well models perform after a small amount of additional training on probes (**inoculation** (Liu et al., 2019))?



**Controls:** Probes should be demonstrably difficult (**strong baselines**);  
Re-training must preserve performance (minimal **inoculation loss**).



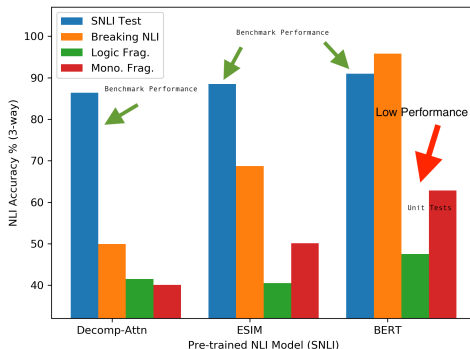
What happens when we do unit testing?  
(NLI  $\rightarrow$  QA)

## Result 1: It's easy to find bugs: zero-shot (NLI)

- ▶ **Approach:** do testing on models trained on benchmark tasks, look at difference in performance; the difference is usually large.

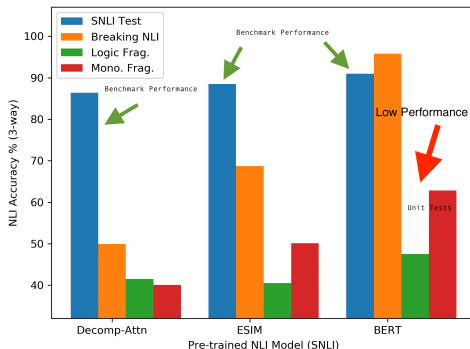
# Result 1: It's easy to find bugs: zero-shot (NLI)

- **Approach:** do testing on models trained on benchmark tasks, look at difference in performance; the difference is usually large.



# Result 1: It's easy to find bugs: zero-shot (NLI)

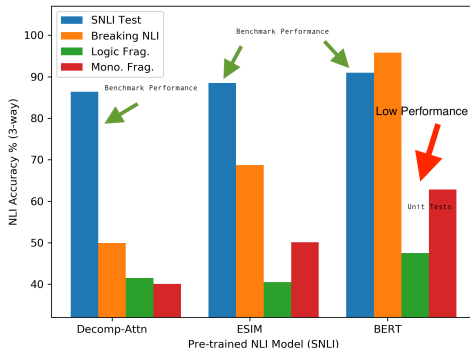
- **Approach:** do testing on models trained on benchmark tasks, look at difference in performance; the difference is usually large.



NLI models trained on standard benchmarks are still lacking in basic linguistic and reasoning abilities.

# Result 1: It's easy to find bugs: zero-shot (NLI)

- **Approach:** do testing on models trained on benchmark tasks, look at difference in performance; the difference is usually large.



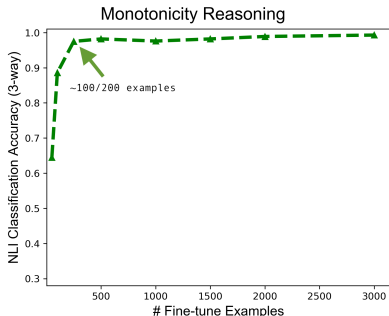
**Caveats:** Do models really not possess the target knowledge, or lack knowledge of format?

## Result 2: Patching bugs can be easy (best NLI models)

- ▶ **Model Inoculation** ([Richardson et al., 2020](#))[AAAI]: Continue training models on small amounts of diagnostic; aim to (quickly) fix model.

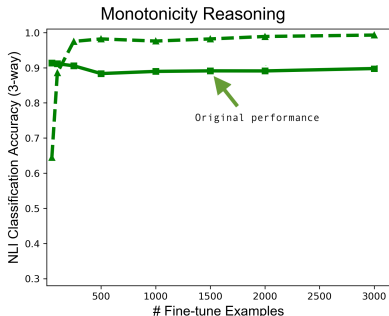
## Result 2: Patching bugs can be easy (best NLI models)

- **Model Inoculation** ([Richardson et al., 2020](#))[AAAI]: Continue training models on small amounts of diagnostic; aim to (quickly) fix model.



## Result 2: Patching bugs can be easy (best NLI models)

- **Model Inoculation** (Richardson et al., 2020)[AAAI]: Continue training models on small amounts of diagnostic; aim to (quickly) fix model.

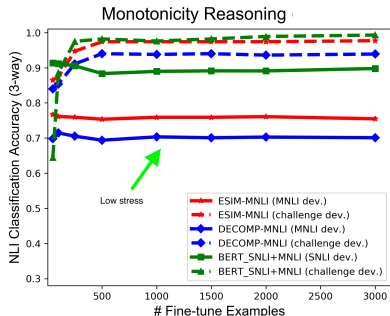


**Assumption:** Ability of model to quickly learn new tasks with minimal effect on original task indicates competence.



## Result 2: Patching bugs can be easy (best NLI models)

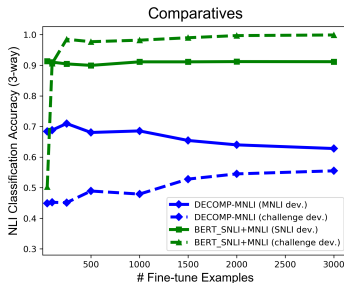
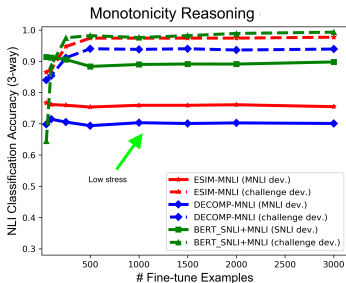
- **Model Inoculation** (Richardson et al., 2020)[AAAI]: Continue training models on small amounts of diagnostic; aim to (quickly) fix model.



**Assumption:** Ability of model to quickly learn new tasks with minimal effect on original task indicates competence.

## Result 2: Patching bugs can be easy (best NLI models)

- **Model Inoculation** (Richardson et al., 2020)[AAAI]: Continue training models on small amounts of diagnostic; aim to (quickly) fix model.



**Assumption:** Ability of model to quickly learn new tasks with minimal effect on original task indicates competence.

## Result 2: Patching bugs can be easy (best NLI models)

- **Model Inoculation** ([Richardson et al., 2020](#))[AAAI]: Continue training models on small amounts of diagnostic; aim to (quickly) fix model.

Category ↓, Model →	BERT Transformer	ESIM	DecAttn
Conditionals	😊	😊	😊
Counting	😊	😊	😞
Quantifiers	😊	😞	😞
Negation	😊	😞	😞
Boolean Coordination	😊	😞	😞
Comparatives	😊	😞	😞

😞 = (bad/mediocre performance + forgetting on **test**), 😊 = (high performance + minimal forgetting on **test**)

## Result 3. Transformers do have lexical knowledge (QA)

- ▶ **Zero-shot**, models do well on *some* categories of knowledge ; far outpace baselines trained on diagnostics.

## Result 3. Transformers do have lexical knowledge (QA)

- **Zero-shot**, models do well on *some* categories of knowledge ; far outpace baselines trained on diagnostics.

Diagnostic performance (QA Accuracy %; random ~ 20%)					
Model	Definitions	Synonymy	Hypernymy	Hyponymy	WordSense
trained LSTM + GloVe	51.8%	55.3%	47.0%	64.2%	53.5%
BERT (zero-shot)	55.7%	60.9%	51.0%	27.0%	42.9%
RoBERTa (zero-shot)	77.1 %	64.2%	71.0%	58.0%	55.1%
Human	91.2%	87.4%	96%	95.5%	95.6%

## Result 3. Transformers do have lexical knowledge (QA)

- **Zero-shot**, models do well on *some* categories of knowledge ; far outpace baselines trained on diagnostics.

Diagnostic performance (QA Accuracy %; random ~ 20%)					
Model	Definitions	Synonymy	Hypernymy	Hyponymy	WordSense
trained LSTM + GloVe	51.8%	55.3%	47.0%	64.2%	53.5%
BERT (zero-shot)	55.7%	60.9%	51.0%	27.0%	42.9%
RoBERTa (zero-shot)	77.1 %	64.2%	71.0%	58.0%	55.1%
Human	91.2%	87.4%	96%	95.5%	95.6%

**Caveats:** Reflect true model knowledge or (non-)familiarity with format? Lower-bound estimate ([Petroni et al., 2019](#)).

## Result 3. Transformers do have lexical knowledge (QA)

- **Inoculation setting:** Models quickly start reaching near human performance .

Diagnostic performance (QA Accuracy %; random ~ 20%)					
Model	Definitions	Synonymy	Hypernymy	Hyponymy	WordSense
BERT (inoculation)	84.1%	79.7%	82.7%	88.0%	79.1%
RoBERTa (inoculation)	89.3 %	81.3%	87.0%	89.4%	85.4%
Human	91.2%	87.4%	96%	95.5%	95.6%

## Result 3. Transformers do have lexical knowledge (QA)

- **Inoculation setting:** Models quickly start reaching near human performance .

Diagnostic performance (QA Accuracy %; random ~ 20%)					
Model	Definitions	Synonymy	Hypernymy	Hyponymy	WordSense
BERT (inoculation)	84.1%	79.7%	82.7%	88.0%	79.1%
RoBERTa (inoculation)	89.3 %	81.3%	87.0%	89.4%	85.4%
Human	91.2%	87.4%	96%	95.5%	95.6%

Giving the model the chance to learn **target format** is important, gives better picture of competence ; minimal loss.



## Result 3. Looking a bit deeper... (QA)

- ▶ The controlled nature of the probes allows for a more granular examination of performance.

## Result 3. Looking a bit deeper... (QA)

- ▶ The controlled nature of the probes allows for a more granular examination of performance.

LSTM Baseline (QA task)

Datasets and # of Hops	hypernyms, k=1	0.27	0.26	0.27	0.23	0.29	0.22	
	hypernyms, k=2	0.29	0.26	0.29	0.31	0.34	0.3	0.33
	hypernyms, k=3	0.3	0.27	0.29	0.27	0.4	0.25	
	hypernyms, k=4	0.29	0.25	0.2	0.25	0.29	0	
	hypernyms, k=5	0.31	0.25	0.2	0.33	0.18		
	hyponyms, k=1	0.33	0.23	0.26	0.23	0.23	0.22	0.25
	hyponyms, k=2	0.29	0.18	0.18	0.2	0.16	0.2	0.2
	hyponyms, k=3	0.39	0.18	0.19	0.15	0.16	0.094	0.17
	hyponyms, k=4	0.091	0	0.21	0	0.17		
	definitions	0.31	0.27	0.31	0.28	0.28	0.27	0.24
	synonyms	0.36	0.22	0.3	0.26	0.21	0.2	0.23

Distractor Types

## Result 3. Looking a bit deeper... (QA)

- ▶ The controlled nature of the probes allows for a more granular examination of performance.

ISA reasoning 3 steps  
moderate distractors

LSTM Baseline (QA task)

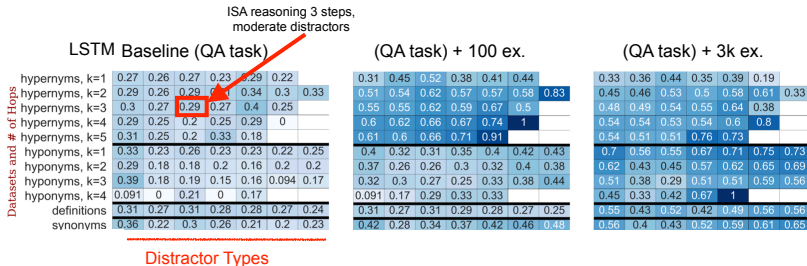
Datasets and # of Hops

hypernyms, k=1	0.27	0.26	0.27	0.23	0.29	0.22	
hypernyms, k=2	0.29	0.26	0.29	0.27	0.34	0.3	0.33
hypernyms, k=3	0.3	0.27	0.29	0.27	0.4	0.25	
hypernyms, k=4	0.29	0.25	0.2	0.25	0.29	0	
hypernyms, k=5	0.31	0.25	0.2	0.33	0.18		
hyponyms, k=1	0.33	0.23	0.26	0.23	0.23	0.22	0.25
hyponyms, k=2	0.29	0.18	0.18	0.2	0.16	0.2	0.2
hyponyms, k=3	0.39	0.18	0.19	0.15	0.16	0.094	0.17
hyponyms, k=4	0.091	0	0.21	0	0.17		
definitions	0.31	0.27	0.31	0.28	0.28	0.27	0.24
synonyms	0.36	0.22	0.3	0.26	0.21	0.2	0.23

Distractor Types

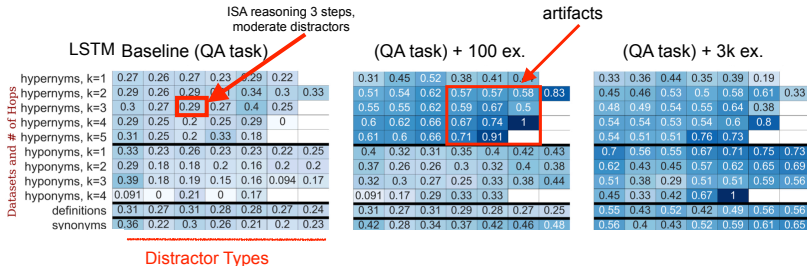
# Result 3. Looking a bit deeper... (QA)

- The controlled nature of the probes allows for a more granular examination of performance.



# Result 3. Looking a bit deeper... (QA)

- The controlled nature of the probes allows for a more granular examination of performance.







# Result 3. Looking a bit deeper... (QA)

- The controlled nature of the probes allows for a more granular examination of performance.

ISA reasoning 3 steps, moderate distractors

**LSTM Baseline (QA task)**

*Datasets and # of Hops*

hypernyms, k=1	0.27	0.26	0.27	0.23	0.29	0.22
hypernyms, k=2	0.29	0.26	0.29	0.27	0.34	0.33
hypernyms, k=3	0.3	0.27	0.29	0.27	0.4	0.25
hypernyms, k=4	0.29	0.25	0.2	0.25	0.29	0
hypernyms, k=5	0.31	0.25	0.2	0.33	0.18	
hyponyms, k=1	0.33	0.23	0.26	0.23	0.23	0.25
hyponyms, k=2	0.29	0.18	0.18	0.2	0.16	0.2
hyponyms, k=3	0.39	0.18	0.19	0.15	0.16	0.094
hyponyms, k=4	0.091	0	0.21	0	0.17	
definitions	0.31	0.27	0.31	0.28	0.28	0.27
synonyms	0.36	0.22	0.3	0.26	0.21	0.2

**(QA task) + 100 ex.**

0.31	0.45	0.52	0.38	0.41	0.44
0.51	0.54	0.62	0.57	0.57	0.58
0.55	0.55	0.62	0.59	0.67	0.5
0.6	0.62	0.66	0.67	0.74	1
0.61	0.6	0.66	0.71	0.91	
0.4	0.32	0.31	0.35	0.4	0.42
0.37	0.26	0.26	0.3	0.32	0.4
0.32	0.3	0.27	0.25	0.33	0.38
0.091	0.17	0.29	0.33	0.33	
0.31	0.27	0.31	0.29	0.28	0.27
0.42	0.28	0.34	0.37	0.42	0.46

**(QA task) + 3k ex.**

0.33	0.36	0.44	0.35	0.39	0.19
0.45	0.46	0.53	0.5	0.58	0.61
0.48	0.49	0.54	0.55	0.64	0.38
0.54	0.54	0.53	0.54	0.6	0.8
0.54	0.51	0.51	0.76	0.73	
0.7	0.56	0.55	0.67	0.71	0.75
0.62	0.43	0.45	0.57	0.62	0.65
0.51	0.38	0.29	0.51	0.51	0.59
0.45	0.33	0.42	0.67	1	
0.55	0.43	0.52	0.42	0.49	0.56
0.56	0.4	0.43	0.52	0.59	0.61

**SOTA Transformer (QA task)**

*Datasets and # of Hops*

hypernyms, k=1	0.76	0.57	0.68	0.48	0.64	0.85
hypernyms, k=2	0.71	0.47	0.58	0.43	0.57	0.7
hypernyms, k=3	0.65	0.38	0.5	0.4	0.54	0.69
hypernyms, k=4	0.61	0.33	0.4	0.33	0.43	0.4
hypernyms, k=5	0.62	0.33	0.46	0.26	0.27	
hyponyms, k=1	0.72	0.47	0.58	0.34	0.46	0.55
hyponyms, k=2	0.59	0.37	0.45	0.25	0.35	0.45
hyponyms, k=3	0.43	0.24	0.3	0.1	0.098	0.19
hyponyms, k=4	0.64	0.15	0.38	0	0.33	
definitions	0.88	0.72	0.8	0.64	0.73	0.76
synonyms	0.82	0.49	0.67	0.63	0.63	0.61

**(QA task) + 100 ex.**

0.83	0.71	0.82	0.61	0.81	0.85
0.8	0.63	0.76	0.61	0.71	0.85
0.77	0.58	0.71	0.59	0.72	0.75
0.79	0.58	0.67	0.56	0.74	0.8
0.79	0.59	0.76	0.41	0.36	
0.89	0.68	0.74	0.64	0.81	0.85
0.77	0.5	0.59	0.43	0.65	0.7
0.65	0.4	0.42	0.34	0.51	0.56
0.55	0.12	0.25	0.33	0.67	
0.94	0.83	0.88	0.7	0.84	0.9
0.85	0.5	0.67	0.73	0.82	0.83

**(QA task) + 3k ex.**

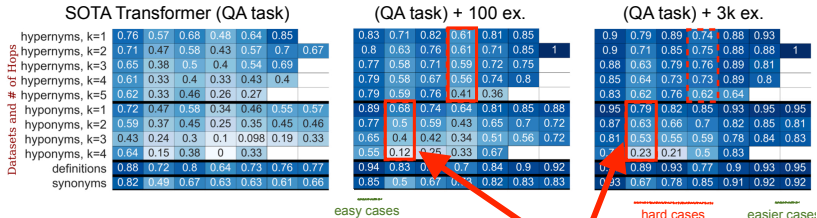
0.9	0.79	0.89	0.74	0.88	0.93
0.9	0.71	0.85	0.75	0.88	0.88
0.88	0.63	0.79	0.76	0.89	0.81
0.85	0.64	0.73	0.73	0.89	0.8
0.83	0.62	0.76	0.62	0.64	
0.95	0.79	0.82	0.85	0.93	0.95
0.87	0.63	0.66	0.7	0.82	0.85
0.81	0.53	0.55	0.59	0.78	0.84
0.73	0.23	0.21	0.5	0.83	
0.97	0.89	0.93	0.77	0.9	0.93
0.93	0.67	0.78	0.85	0.91	0.92

Can nudge models to bring out knowledge with small set of examples, cheap way to inject knowledge into transformers, can be used as KBs.



# Result 3. Looking a bit deeper... (QA)

- ▶ The controlled nature of the probes allows for a more granular examination of performance.

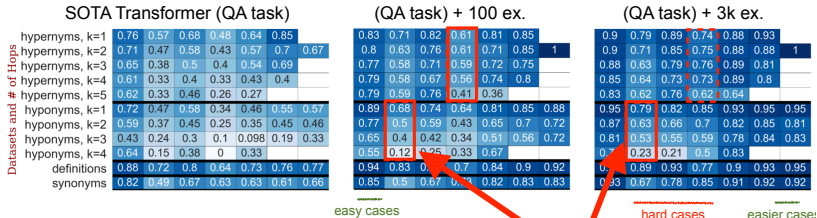


Several inference steps

Model does show sensitivity to reasoning complexity, weak areas.

# Result 3. Looking a bit deeper... (QA)

- ▶ The controlled nature of the probes allows for a more granular examination of performance.

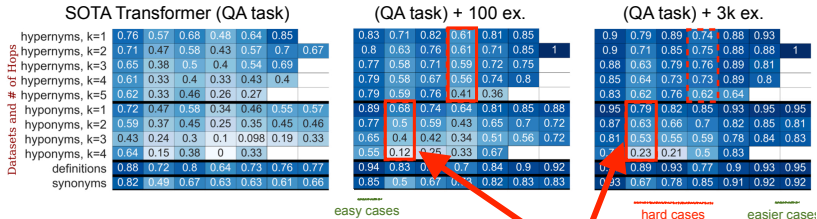


Model does show sensitivity to reasoning complexity; weak areas.

1. have knowledge across many concepts; ✓
2. be robust to perturbations ✓ / ?
3. and varying levels of reasoning complexity ?

# Result 3. Looking a bit deeper... (QA)

- ▶ The controlled nature of the probes allows for a more granular examination of performance.



Several inference steps

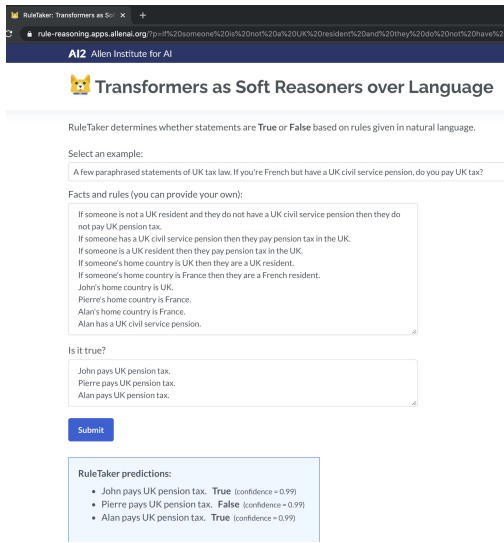
Model does show **sensitivity to reasoning complexity**; weak areas.

1. have knowledge across *many concepts*; ✓
2. be robust to *perturbations* ✓ / ?
3. and varying levels of reasoning *complexity* ?

**Probing is difficult!** Definitive proof of model knowledge is difficult.

# Transformers for Formalized (Deductive) Reasoning

- ▶ Achieve high accuracy on reasoning tasks (99% accuracy); Ability to generalize to human authored theories, reasoning depths.



The screenshot shows a web browser window with the address bar displaying "rule-reasoning.apps.allenai.org". The page header includes the Allen Institute for AI logo and the title "Transformers as Soft Reasoners over Language". The main content area explains that RuleTaker determines the truth of statements based on rules. It provides an example scenario involving UK tax law and pension rules. The user is asked to select an example and provide facts and rules. The example facts and rules are: "If someone is not a UK resident and they do not have a UK civil service pension then they do not pay UK pension tax.", "If someone has a UK civil service pension then they pay pension tax in the UK.", "If someone is a UK resident then they pay pension tax in the UK.", "If someone's home country is UK then they are a UK resident.", "If someone's home country is France then they are a French resident.", "John's home country is UK.", "Pierre's home country is France.", "Alan's home country is France.", "Alan has a UK civil service pension." The user is asked "Is it true?" and provides the statements: "John pays UK pension tax.", "Pierre pays UK pension tax.", "Alan pays UK pension tax." The "Submit" button is highlighted. The "RuleTaker predictions:" section shows the results: "John pays UK pension tax. True (confidence = 0.99)", "Pierre pays UK pension tax. False (confidence = 0.99)", and "Alan pays UK pension tax. True (confidence = 0.99)".

RuleTaker: Transformers as Soft Reasoners over Language

rule-reasoning.apps.allenai.org

AI2 Allen Institute for AI

## Transformers as Soft Reasoners over Language

RuleTaker determines whether statements are **True** or **False** based on rules given in natural language.

Select an example:

A few paraphrased statements of UK tax law. If you're French but have a UK civil service pension, do you pay UK tax?

Facts and rules (you can provide your own):

If someone is not a UK resident and they do not have a UK civil service pension then they do not pay UK pension tax.  
If someone has a UK civil service pension then they pay pension tax in the UK.  
If someone is a UK resident then they pay pension tax in the UK.  
If someone's home country is UK then they are a UK resident.  
If someone's home country is France then they are a French resident.  
John's home country is UK.  
Pierre's home country is France.  
Alan's home country is France.  
Alan has a UK civil service pension.

Is it true?

John pays UK pension tax.  
Pierre pays UK pension tax.  
Alan pays UK pension tax.

Submit

RuleTaker predictions:

- John pays UK pension tax. **True** (confidence = 0.99)
- Pierre pays UK pension tax. **False** (confidence = 0.99)
- Alan pays UK pension tax. **True** (confidence = 0.99)

</Results>

# Conclusions

- ▶ Probing using symbolic models, tasks with 1. wide range of concepts 2. systematic perturbations ; 3. variable complexity.
  - ▶ Useful for better understanding models, supplement to existing NLU research; **few-shot learning** for model fixing.
  - ▶ Inherently collaborative enterprise, need help!
- ▶ several diagnostic tasks for QA, NLI; **extending to other tasks**, behavioral testing + **interventions** (manipulations of network states) ([Geiger et al., 2020](#))[BlackBoxNLP].

# Conclusions

- ▶ Probing using symbolic models, tasks with 1. wide range of concepts 2. systematic perturbations ; 3. variable complexity.
  - ▶ Useful for better understanding models, supplement to existing NLU research; **few-shot learning** for model fixing.
  - ▶ Inherently collaborative enterprise, need help!
- ▶ several diagnostic tasks for QA, NLI; **extending to other tasks**, behavioral testing + **interventions** (manipulations of network states) ([Geiger et al., 2020](#))[BlackBoxNLP].

**Tooling:** allowing non-experts to author their own datasets, *democratize* the dataset construction process.

Thank you.



# References I

- Boratto, M., Padigela, H., Mikkilineni, D., Yuvraj, P., Das, R., McCallum, A., Chang, M., Fokoue-Nkoutche, A., Kapanipathi, P., Mattei, N., et al. (2018). A systematic classification of knowledge, reasoning, and context within the arc dataset. *arXiv preprint arXiv:1806.00358*.
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. (2018). Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- Clark, P., Tafjord, O., and Richardson, K. (2020). Transformers as soft reasoners over language. *Proceedings of IJCAI*.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Geiger, A., Richardson, K., and Potts, C. (2020). Modular representation underlies systematic generalization in neural natural language inference models. *arXiv preprint arXiv:2004.14623*.
- Glockner, M., Shwartz, V., and Goldberg, Y. (2018). Breaking NLI Systems with Sentences that Require Simple Lexical Inferences. *arXiv preprint arXiv:1805.02266*.
- Gururangan, S., Swayamdipta, S., Levy, O., Schwartz, R., Bowman, S., and Smith, N. A. (2018). Annotation Artifacts in Natural Language Inference Data. In *Proceedings of NAACL: HLT*.

## References II

- Khot, T., Khashabi, D., Richardson, K., Clark, P., and Sabharwal, A. (2020). Text Modular Networks: Learning to Decompose Tasks in the Language of Existing Models. *arXiv preprint arXiv:2009.00751*.
- Liu, N. F., Schwartz, R., and Smith, N. A. (2019). Inoculation by Fine-Tuning: A Method for Analyzing Challenge Datasets. *arXiv preprint arXiv:1904.02668*.
- Mihaylov, T., Clark, P., Khot, T., and Sabharwal, A. (2018). Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*.
- Montague, R. (1973). The proper treatment of quantification in ordinary english. In *Approaches to natural language*, pages 221–242. Springer.
- Naik, A., Ravichander, A., Sadeh, N., Rose, C., and Neubig, G. (2018). Stress test evaluation for natural language inference. *arXiv preprint arXiv:1806.00692*.
- Petroni, F., Rocktäschel, T., Lewis, P., Bakhtin, A., Wu, Y., Miller, A. H., and Riedel, S. (2019). Language Models as Knowledge Bases? *arXiv preprint arXiv:1909.01066*.
- Pratt-Hartmann, I. (2004). Fragments of language. *Journal of Logic, Language and Information*, 13(2):207–223.
- Pratt-Hartmann, I. (2015). Semantic complexity in natural language. *The Handbook of Contemporary Semantic Theory*, page 429.
- Ribeiro, M. T., Wu, T., Guestrin, C., and Singh, S. (2020). Beyond accuracy: Behavioral testing of nlp models with checklist. *Proceedings of ACL*.

## References III

- Richardson, K., Hu, H., Moss, L. S., and Sabharwal, A. (2020). Probing Natural Language Inference Models through Semantic Fragments. In *AAAI*, pages 8713–8721.
- Richardson, K. and Sabharwal, A. (2020). What Does My QA Model Know? Devising Controlled Probes using Expert Knowledge. *to appear in TACL*.
- Rozen, O., Shwartz, V., Aharoni, R., and Dagan, I. (2019). Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets. *arXiv preprint arXiv:1910.09302*.
- van Benthem, J. (1986). *Essays in Logical Semantics*, volume 29 of *Studies in Linguistics and Philosophy*. D. Reidel Publishing Co., Dordrecht.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.