

# Do NLP Models Know Numbers? Probing Numeracy in Embeddings

Empirical Methods in Natural Language Processing(EMNLP), 2019

Eric Wallace, Yizhong Wang, Sujian Li  
Sameer Singh, Matt Gardner

**Soojung Kim**  
April 20, 2020

# Before start

왜 이 논문을 선택하게 됐는지



**강산농원 어우티 1.5g x 20개입**

**최저 15,800원** 판매처 4

식품 > 음료 > 차류 > 기타차

리뷰 ★★★★★ 5,587 · 등록일 2020.04. · 찜하기 9 · 정보 수정요청

쇼핑몰별 최저가

티트리트 	↓ 15,800
G마켓	31,640
CJmall 	35,150

- 1개 15,800원 | 1개

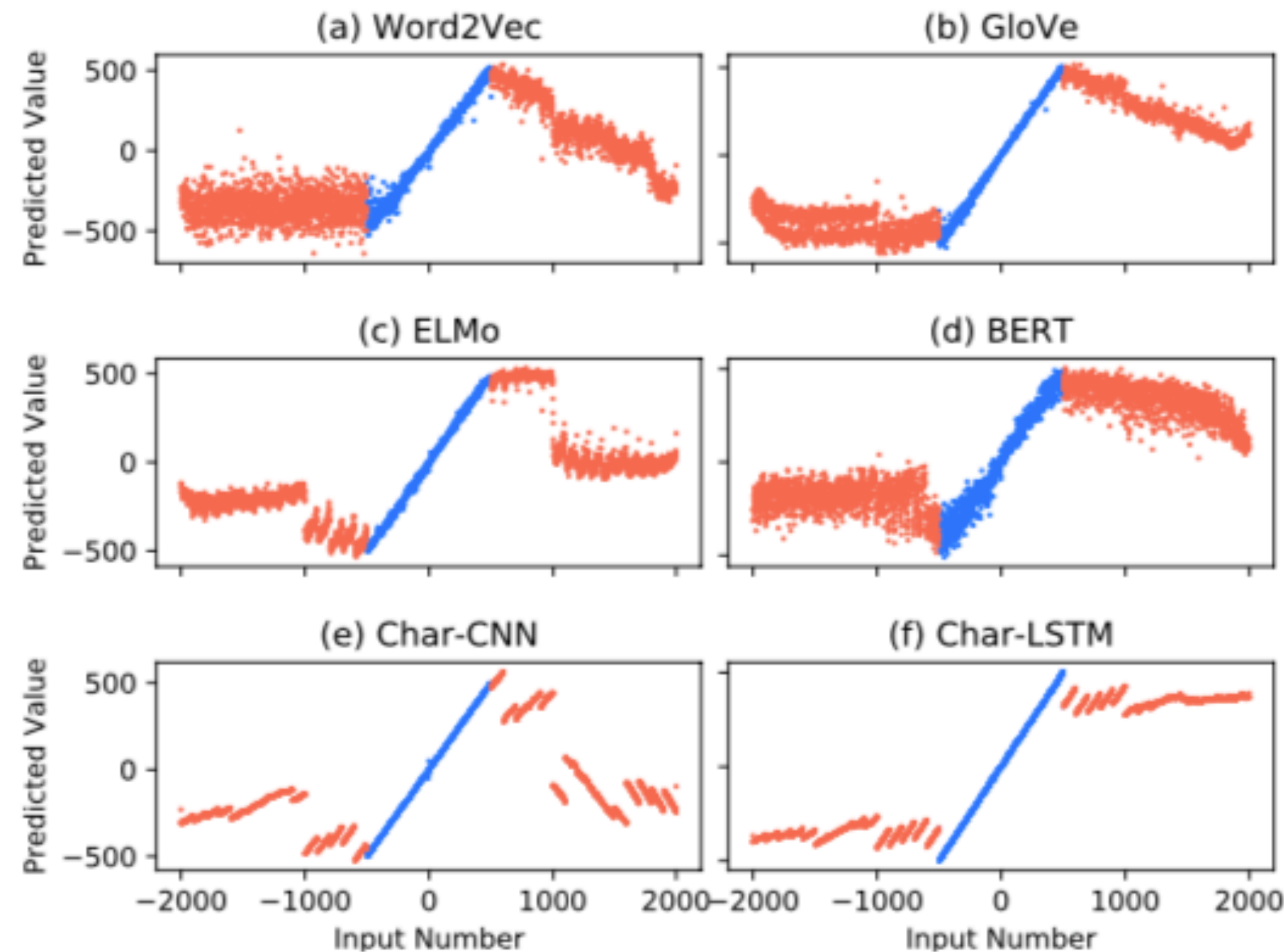
- 2개 29,700원 | 3개

# Abstract

NLP model treat numbers in the same way as other tokens. e.g., distributed vector

Q. Is this a good representation for numerical tasks?

A: Yes! Pre-trained vectors(Glove, ELMo, BERT) know numbers



# Introduction

---

Numerical reasoning ability:  
the ability to understand and work with numbers in digit or word form

String "23" is bigger than "twenty-two".

# Introduction

Numerical reasoning ability:

Extracting the list of field goals and computing that list's maximum

...JaMarcus Russell completed a 91-yard touchdown pass to rookie wide receiver Chaz Schilens. The Texans would respond with fullback Vonta Leach getting a 1-yard touchdown run, yet the Raiders would answer with kicker Sebastian Janikowski getting a 33-yard and a 21-yard field goal. Houston would tie the game in the second quarter with kicker Kris Brown getting a 53-yard and a 24-yard field goal. Oakland would take the lead in the third quarter with wide receiver Johnnie Lee Higgins catching a 29-yard touchdown pass from Russell, followed up by an 80-yard punt return for a touchdown.

Q: How many yards was the longest field goal? A: 53

Q: How long was the shortest touchdown pass? A: 29-yard

Q: Who caught the longest touchdown? A: Chaz Schilens

*if and how NLP models learn to reason numerically  
over paragraphs with only question-answer supervision*

1. Analyzing the sota NAQANet model
2. Probe token embedding methods

<https://demo.allennlp.org/reading-comprehension/MTczNDYlMw==>



# Analyzing the sota NAQANet model

Testing QA models on questions that evaluate numerical reasoning taken DROP dataset (e.g., sorting, comparing, or summing numbers.)

...JaMarcus Russell completed a 91-yard touchdown pass to rookie wide receiver Chaz Schilens. The Texans would respond with fullback Vonta Leach getting a 1-yard touchdown run, yet the Raiders would answer with kicker Sebastian Janikowski getting a 33-yard and a 21-yard field goal. Houston would tie the game in the second quarter with kicker Kris Brown getting a 53-yard and a 24-yard field goal. Oakland would take the lead in the third quarter with wide receiver Johnnie Lee Higgins catching a 29-yard touchdown pass from Russell, followed up by an 80-yard punt return for a touchdown.

Q: How many yards was the longest field goal? A: 53

Q: How long was the shortest touchdown pass? A: 29-yard

Q: Who caught the longest touchdown? A: Chaz Schilens

- DROP dataset
- NAQANet (sota)

sorting

comparing



# DROP dataset: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs

Reading comprehension is the task of answering questions about a passage of text to show that **the system understands the passage.**

Reasoning	Passage (some parts shortened)	Question	Answer	BiDAF
Subtraction (28.8%)	That year, his <b>Untitled (1981)</b> , a painting of a haloed, black-headed man with a bright red skeletal body, depicted amid the artists signature scrawls, was <b>sold by Robert Lehrman for \$16.3 million, well above its \$12 million high estimate.</b>	How many more dol-lars was the Untitled (1981) painting sold for than the 12 million dollar estimation?	4300000	\$16.3 million
Comparison (18.2%)	In <b>1517, the seventeen-year-old King sailed to Castile.</b> There, his Flemish court . . . . <b>In May 1518, Charles traveled to Barcelona in Aragon.</b>	Where did Charles travel to first, Castile or Barcelona?	Castile	Aragon
Selection (19.4%)	In 1970, to commemorate the 100th anniversary of the founding of Baldwin City, <b>Baker University professor and playwright Don Mueller and Phyllis E. Braun, Business Manager, produced a musical play entitled The Ballad Of Black Jack</b> to tell the story of the events that led up to the battle.	Who was the Uni-versity professor that helped produce The Ballad Of Black Jack, Ivan Boyd or Don Mueller?	Don Mueller	Baker
Addition (11.7%)	Before the UNPROFOR fully deployed, the HV clashed with an armed force of the RSK in the village of Nos Kalik, located in a pink zone near Šibenik, and captured the village at 4:45 p.m. on <b>2 March 1992.</b> The JNA formed a battlegroup to counterattack the <b>next day.</b>	What date did the JNA form a battlegroup to counterattack after the village of Nos Kalik was captured?	3 March 1992	2 March 1992
Count (16.5%) and Sort (11.7%)	Denver would retake the lead with kicker <b>Matt Prater nailing a 43-yard field goal</b> , yet Carolina answered as kicker <b>John Kasay ties the game with a 39-yard field goal.</b> . . . Carolina closed out the half with <b>Kasay nail-ing a 44-yard field goal.</b> . . . In the fourth quarter, Car-olina sealed the win with <b>Kasay’s 42-yard field goal.</b>	Which kicker kicked the most field goals?	John Kasay	Matt Prater
Coreference Resolution (3.7%)	<b>James Douglas</b> was the second son of Sir George Dou-glas of Pittendreich, and Elizabeth Douglas, daughter David Douglas of Pittendreich. Before <b>1543 he mar-ried Elizabeth</b> , daughter of James Douglas, 3rd Earl of Morton. <b>In 1553 James Douglas succeeded to the title and estates of his father-in-law.</b>	How many years af-ter he married Eliza-beth did James Dou-glas succeed to the ti-tle and estates of his father-in-law?	10	1553
Other Arithmetic (3.2%)	Although the movement initially gathered some <b>60,000 adherents</b> , the subsequent establishment of the Bulgar-ian Exarchate <b>reduced their number by some 75%.</b>	How many adherents were left after the es-tablishment of the Bul-garian Exarchate?	15000	60,000
Set of spans (6.0%)	According to some sources 363 civilians were killed in <b>Kavadarci</b> , 230 in <b>Negotino</b> and 40 in <b>Vatasha.</b>	What were the 3 vil-lages that people were killed in?	Kavadarci, Negotino, Vatasha	Negotino and 40 in Vatasha
Other (6.8%)	This <b>Annual Financial Report</b> is our principal financial statement of accountability. The <b>AFR gives a compre-hensive view</b> of the Department’s financial activities ...	What does AFR stand for?	Annual Financial Report	one of the Big Four audit firms

# DROP dataset: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs

focus on Comparative and Superlative questions

## Comparative questions

- quantities or events "larger", "smaller", "longer", ...
- "either-or" relations

Question Type	Example	Reasoning Required
Comparative (Binary)	Which country is a bigger exporter, Brazil or Uruguay?	Binary Comparison
Comparative (Non-binary)	Which player had a touchdown longer than 20 yards?	Greater Than

## Superlative questions

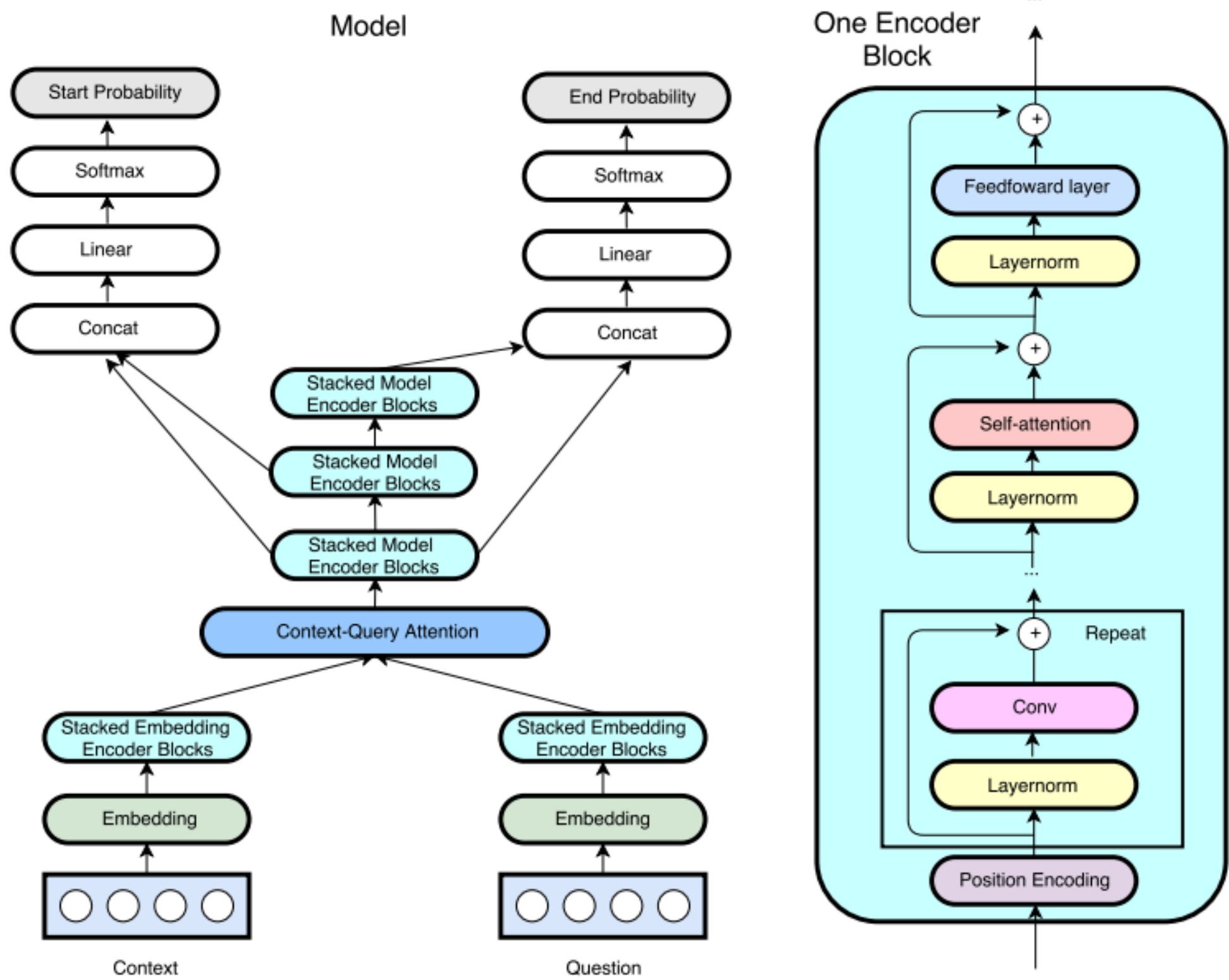
- "shortest", "longest", "biggest"
- finding the maximum or minimum of a list
- argmax operation

Superlative (Number)	How many yards was the shortest field goal?	List Minimum
Superlative (Span)	Who kicked the longest field goal?	Argmax



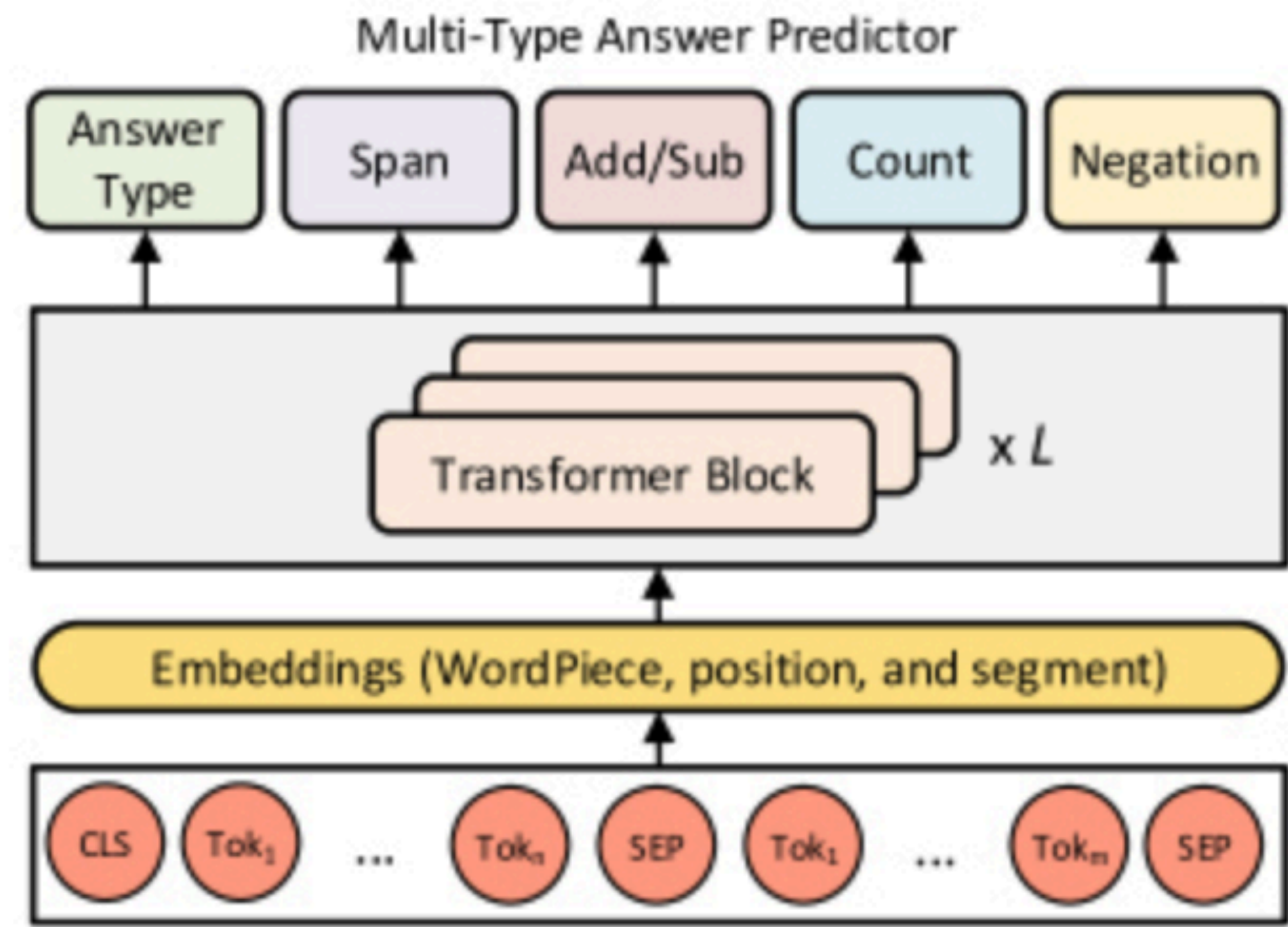
# NAQANet model

QANet + four output branches (passage span, question span, count, add/sub)



QANet

<https://arxiv.org/pdf/1804.09541.pdf>



Multi type answer predictor

<https://deepai.org/publication/a-multi-type-multi-span-network-for-reading-comprehension-that-requires-discrete-reasoning>

# Analyzing the sota NAQANet model

NAQANet achieving only 49 F1 on the entire validation set  
It scores 89 F1 on numerical comparsion questions.

Question Type	Count	EM	F1
Human (Test Set)	9622	92.4	96.0
Full Validation	9536	46.2	49.2
Number Answers	5842	44.3	44.4
Comparative	704	73.6	76.4
Binary (either-or)	477	86.0	89.0
Non-binary	227	47.6	49.8
Superlative Questions	861	64.6	67.7
Number Answers	475	68.8	69.2
Span Answers	380	59.7	66.3

# Stress Testing NAQANet's Numeracy

NAQANet has a strong understanding of numeracy for numbers in the training range, but the model can **fail** to extrapolate to other values

training range [0, 100]

Stress Test Dataset	All Questions		Superlative	
	F1	$\Delta$	F1	$\Delta$
Original Validation Set	49.2	-	67.7	-
Add [1, 20]	47.7	-1.5	64.1	-3.6
Add [21, 100]	41.4	-7.8	40.4	-27.3
Multiply [2, 10]	41.1	-8.1	39.3	-28.4
Multiply [11, 100]	38.8	-10.4	32.0	-35.7
Digits to Words [0, 20]	45.5	-3.7	63.8	-3.9
Digits to Words [21, 100]	41.9	-7.3	46.1	-21.6

\* Digits to Words: "75"  $\rightarrow$  "seventy-five"

1. Analyzing the sota NAQANet model
- 2. Probe token embedding methods**

<https://demo.allennlp.org/reading-comprehension/MTczNDYlMw==>



# Probe token embedding methods(e.g., BERT, GloVe)

To understand how models can capture numerical reasoning, we use three synthetic tasks to probe the numeracy of token embeddings.

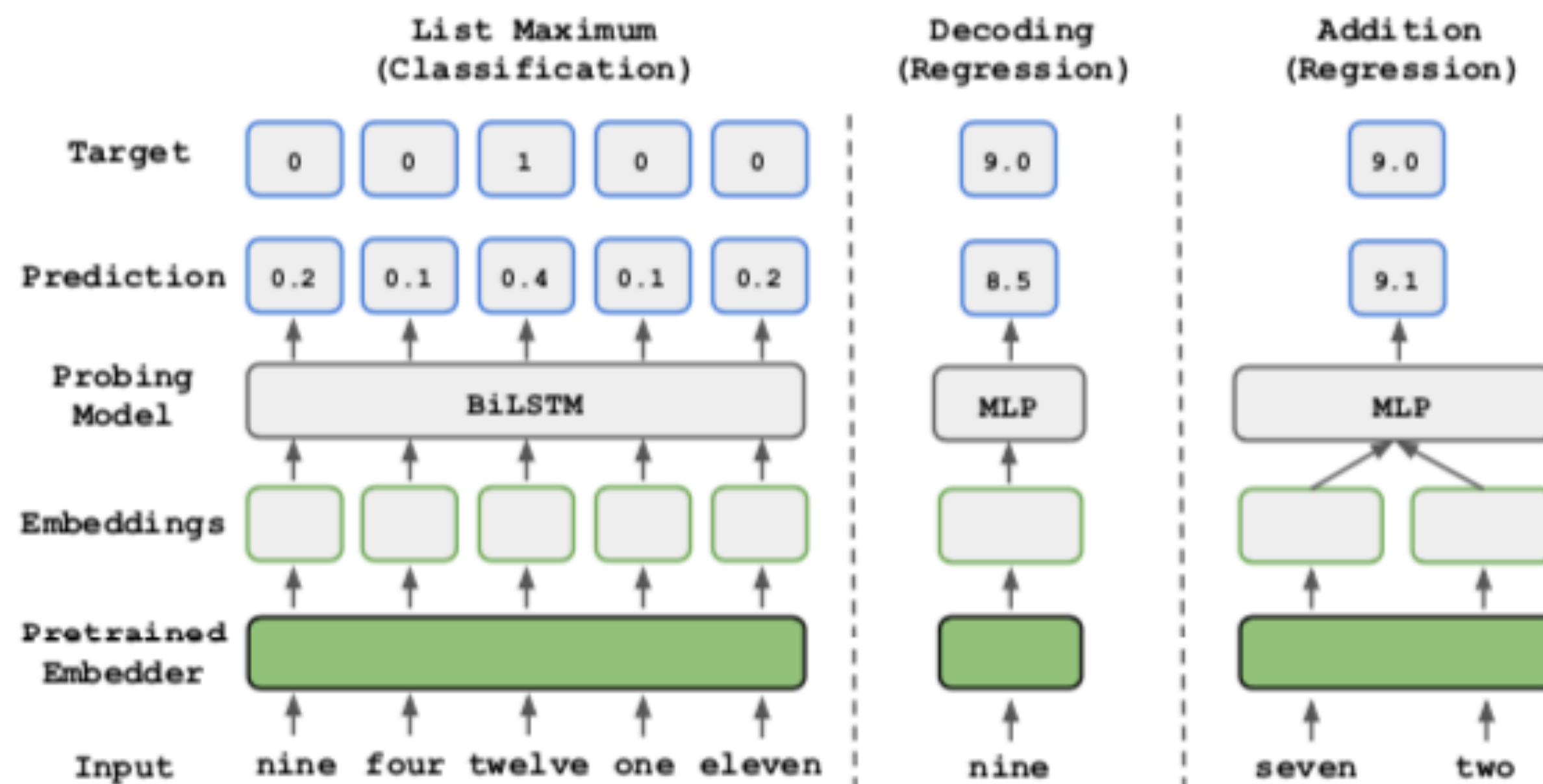


Figure 3: Our probing setup. We pass numbers through a pre-trained embedder (e.g., BERT, GloVe) and train a probing model to solve numerical tasks such as finding a list's maximum, decoding a number, or adding two numbers. If the probing model generalizes to held-out numbers, the pre-trained embeddings must contain numerical information. We provide numbers as either words (shown here), digits ("9"), floats ("9.1"), or negatives ("-9").



## Results: Embeddings capture numeracy

word2vec and GloVe:

continuous bag of words objective can teach fine-grained number magnitude.

<b>Interpolation</b> <i>Integer Range</i>	<b>List Maximum (5-classes)</b>			<b>Decoding (RMSE)</b>			<b>Addition (RMSE)</b>		
	[0,99]	[0,999]	[0,9999]	[0,99]	[0,999]	[0,9999]	[0,99]	[0,999]	[0,9999]
Random Vectors	0.16	0.23	0.21	29.86	292.88	2882.62	42.03	410.33	4389.39
Untrained CNN	0.97	0.87	0.84	2.64	9.67	44.40	1.41	14.43	69.14
Untrained LSTM	0.70	0.66	0.55	7.61	46.5	210.34	5.11	45.69	510.19
Value Embedding	<b>0.99</b>	0.88	0.68	<b>1.20</b>	11.23	275.50	<b>0.30</b>	15.98	654.33
<i>Pre-trained</i>									
Word2Vec	0.90	0.78	0.71	2.34	18.77	333.47	0.75	21.23	210.07
GloVe	0.90	0.78	0.72	2.23	13.77	174.21	0.80	16.51	180.31
ELMo	0.98	0.88	0.76	2.35	13.48	62.20	0.94	15.50	45.71
BERT	0.95	0.62	0.52	3.21	29.00	431.78	4.56	67.81	454.78
<i>Learned</i>									
Char-CNN	0.97	<b>0.93</b>	<b>0.88</b>	2.50	<b>4.92</b>	<b>11.57</b>	1.19	<b>7.75</b>	<b>15.09</b>
Char-LSTM	0.98	0.92	0.76	2.55	8.65	18.33	1.21	15.11	25.37
<i>DROP-trained</i>									
NAQANet	0.91	0.81	0.72	2.99	14.19	62.17	1.11	11.33	90.01
- GloVe	0.88	0.90	0.82	2.87	5.34	35.39	1.45	9.91	60.70

# Results: Embeddings capture numeracy

Character-level methods

strength of the char-level convolutions seems to lie in the architectural prior

<b>Interpolation</b> <i>Integer Range</i>	<b>List Maximum (5-classes)</b>			<b>Decoding (RMSE)</b>			<b>Addition (RMSE)</b>		
	[0,99]	[0,999]	[0,9999]	[0,99]	[0,999]	[0,9999]	[0,99]	[0,999]	[0,9999]
Random Vectors	0.16	0.23	0.21	29.86	292.88	2882.62	42.03	410.33	4389.39
Untrained CNN	0.97	0.87	0.84	2.64	9.67	44.40	1.41	14.43	69.14
Untrained LSTM	0.70	0.66	0.55	7.61	46.5	210.34	5.11	45.69	510.19
Value Embedding	<b>0.99</b>	0.88	0.68	<b>1.20</b>	11.23	275.50	<b>0.30</b>	15.98	654.33
<i>Pre-trained</i>									
Word2Vec	0.90	0.78	0.71	2.34	18.77	333.47	0.75	21.23	210.07
GloVe	0.90	0.78	0.72	2.23	13.77	174.21	0.80	16.51	180.31
ELMo	0.98	0.88	0.76	2.35	13.48	62.20	0.94	15.50	45.71
BERT	0.95	0.62	0.52	3.21	29.00	431.78	4.56	67.81	454.78
<i>Learned</i>									
Char-CNN	0.97	<b>0.93</b>	<b>0.88</b>	2.50	<b>4.92</b>	<b>11.57</b>	1.19	<b>7.75</b>	<b>15.09</b>
Char-LSTM	0.98	0.92	0.76	2.55	8.65	18.33	1.21	15.11	25.37
<i>DROP-trained</i>									
NAQANet	0.91	0.81	0.72	2.99	14.19	62.17	1.11	11.33	90.01
- GloVe	0.88	0.90	0.82	2.87	5.34	35.39	1.45	9.91	60.70



# Results: Embeddings capture numeracy

Sub-word models:

poor method to encode digits: similar two numbers can have very different sub-word divisions

Interpolation <i>Integer Range</i>	List Maximum (5-classes)			Decoding (RMSE)			Addition (RMSE)		
	[0,99]	[0,999]	[0,9999]	[0,99]	[0,999]	[0,9999]	[0,99]	[0,999]	[0,9999]
Random Vectors	0.16	0.23	0.21	29.86	292.88	2882.62	42.03	410.33	4389.39
Untrained CNN	0.97	0.87	0.84	2.64	9.67	44.40	1.41	14.43	69.14
Untrained LSTM	0.70	0.66	0.55	7.61	46.5	210.34	5.11	45.69	510.19
Value Embedding	<b>0.99</b>	0.88	0.68	<b>1.20</b>	11.23	275.50	<b>0.30</b>	15.98	654.33
<i>Pre-trained</i>									
Word2Vec	0.90	0.78	0.71	2.34	18.77	333.47	0.75	21.23	210.07
GloVe	0.90	0.78	0.72	2.23	13.77	174.21	0.80	16.51	180.31
ELMo	0.98	0.88	0.76	2.35	13.48	62.20	0.94	15.50	45.71
BERT	0.95	0.62	0.52	3.21	29.00	431.78	4.56	67.81	454.78
<i>Learned</i>									
Char-CNN	0.97	<b>0.93</b>	<b>0.88</b>	2.50	<b>4.92</b>	<b>11.57</b>	1.19	<b>7.75</b>	<b>15.09</b>
Char-LSTM	0.98	0.92	0.76	2.55	8.65	18.33	1.21	15.11	25.37
<i>DROP-trained</i>									
NAQANet	0.91	0.81	0.72	2.99	14.19	62.17	1.11	11.33	90.01
- GloVe	0.88	0.90	0.82	2.87	5.34	35.39	1.45	9.91	60.70

## Results: Embeddings capture numeracy

A linear subspace exists:  
for small ranges on the decoding task, this capture number magnitude.

<b>Interpolation</b> <i>Integer Range</i>	<b>Decoding (RMSE)</b>		
	[0,50]	[-50,50]	[0,999]
Random	13.86	29.46	275.41
Word2Vec	4.15	8.93	29.04
GloVe	3.21	5.76	23.27
ELMo	1.20	2.89	21.53
BERT	3.23	7.86	64.42

Table 8: *Number Decoding interpolation accuracy with linear regression.* Linear regression is competitive to the fully connected probe for smaller numbers.

## Results: Embeddings capture numeracy

Extrapolation on list maximum.  
trained on the integer range [0, 150] and evaluate on integers from the Test Range

Extrapolation <i>Test Range</i>	List Maximum (5-classes)		
	[151,160]	[151,180]	[151,200]
Rand. Vectors	0.17	0.22	0.15
Untrained CNN	0.80	0.47	0.41
<i>Pre-trained</i>			
Word2Vec	0.14	0.16	0.11
GloVe	0.19	0.17	0.21
ELMo	0.65	0.57	0.38
BERT	0.35	0.11	0.14
<i>Learned</i>			
Char-CNN	0.81	0.75	0.73
Char-LSTM	<b>0.88</b>	<b>0.84</b>	<b>0.82</b>
<i>DROP</i>			
NAQANet	0.31	0.29	0.25
- GloVe	0.58	0.53	0.48



## Results: Embeddings capture numeracy

<b>Interpolation</b> <i>Float Range</i>	<b>List Maximum (5-classes)</b>	
	[0.0,99.9]	[0.0,999.9]
Rand. Vectors	0.18 $\pm$ 0.03	0.21 $\pm$ 0.04
ELMo	0.91 $\pm$ 0.03	0.59 $\pm$ 0.01
BERT	0.82 $\pm$ 0.05	0.51 $\pm$ 0.04
Char-CNN	0.87 $\pm$ 0.04	0.75 $\pm$ 0.03
Char-LSTM	0.81 $\pm$ 0.05	0.69 $\pm$ 0.02

Table 5: *Interpolation with floats (e.g., “18.1”) for list maximum.* Pre-trained embeddings capture numeracy for float values. The probing model is trained on a randomly shuffled 80% of the *Float Range* and tested on the remaining 20%. See the text for details on selecting decimal values. We show the mean alongside the standard deviation over 5 different random shuffles.

<b>Interpolation</b> <i>Integer Range</i>	<b>List Maximum (5-classes)</b>
	[-50,50]
Rand. Vectors	0.23 $\pm$ 0.12
Word2Vec	0.89 $\pm$ 0.02
GloVe	0.89 $\pm$ 0.03
ELMo	0.96 $\pm$ 0.01
BERT	0.94 $\pm$ 0.02
Char-CNN	0.95 $\pm$ 0.07
Char-LSTM	0.97 $\pm$ 0.02

Table 6: *Interpolation with negatives (e.g., “-18”) on list maximum.* Pre-trained embeddings capture numeracy for negative values.

# Augmenting Data to aid Extrapolation

	Superlative		Comparative		All Validation	
	Original	Bigger	Original	Bigger	Original	Bigger
NAQANet	64.5 / 67.7	30.0 / 32.2	73.6 / 76.4	70.3 / 73.0	<b>46.2 / 49.2</b>	38.7 / 41.4
+ Data Augmentation	<b>67.6 / 70.9</b>	<b>59.2 / 62.4</b>	<b>76.0 / 77.7</b>	<b>75.0 / 76.8</b>	46.1 / <b>49.3</b>	<b>42.8 / 45.8</b>

Table 11: Data augmentation improves NAQANet’s interpolation and extrapolation results. We created the *Bigger* version of DROP by multiplying numbers in the passage by a random integer from [11, 20] and then adding a random integer from [21, 40]. Scores are shown in EM / F1 format.



# Conclusion

pre-trained token representations naturally encode numeracy.  
But, model fails to extrapolate to numbers outside the training range

### Reading Comprehension

Reading comprehension is the task of answering questions about a passage of text to show that the system understands the passage.

[Demo](#)[Usage](#)

Enter text or

How many yards was the second longest passing touchdown?

Passage

Hoping to rebound from their loss to the Patriots, the Raiders stayed at home for a Week 16 duel with the Houston Texans. Oakland would get the early lead in the first quarter as quarterback JaMarcus Russell completed a 20-yard touchdown pass to rookie wide receiver Chaz Schilens. The Texans would respond with fullback Vonta Leach getting a 1-yard touchdown run, yet the Raiders would answer with kicker Sebastian Janikowski getting a 33-yard and a 30-yard field goal. Houston would tie the game in the second quarter with kicker Kris Brown getting a 53-yard and a 24-yard field goal. Oakland

Question

How many yards was the second longest passing touchdown?

Model

☐ ELMo-BiDAF (trained on SQuAD)  
☐ BiDAF (trained on SQuAD)  
☒ NAQANet (trained on DROP)  
☐ NMN (trained on DROP)

Run >

