

Designing and Interpreting Probes with Control Tasks



John Hewitt

Percy Liang

Overview

An emerging body of NLP work asks
“Does my neural network implicitly learn task Y?”

Overview

parts-of-speech
syntax
semantics

An emerging body of NLP work asks

“Does my neural network implicitly learn task Y?”

Overview

parts-of-speech
syntax
semantics

An emerging body of NLP work asks
“Does my neural network implicitly learn task Y?”

If a neural network hasn't learned some task,
our methods shouldn't tell us it has.

Overview

parts-of-speech
syntax
semantics

An emerging body of NLP work asks
“Does my neural network implicitly learn task Y?”

If a neural network hasn't learned some task,
our methods shouldn't tell us it has.

(Avoid false positives -- This is hard!)

Probing: supervised analysis of representations

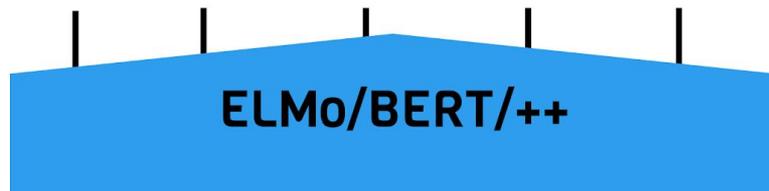
Probing: supervised analysis of representations

Does my network make task (e.g., part-of-speech) labels accessible?

Probing: supervised analysis of representations

Does my network make task (e.g., part-of-speech) labels accessible?

The chef made five pizzas

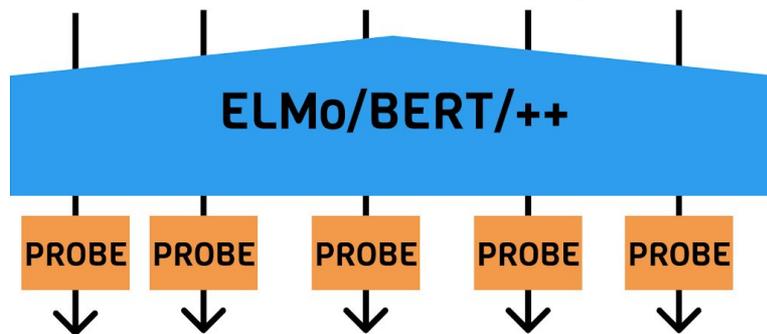


Probing: supervised analysis of representations

Does my network make task (e.g., part-of-speech) labels accessible?

Choose a function family to decode the task. (e.g., linear)

The chef made five pizzas

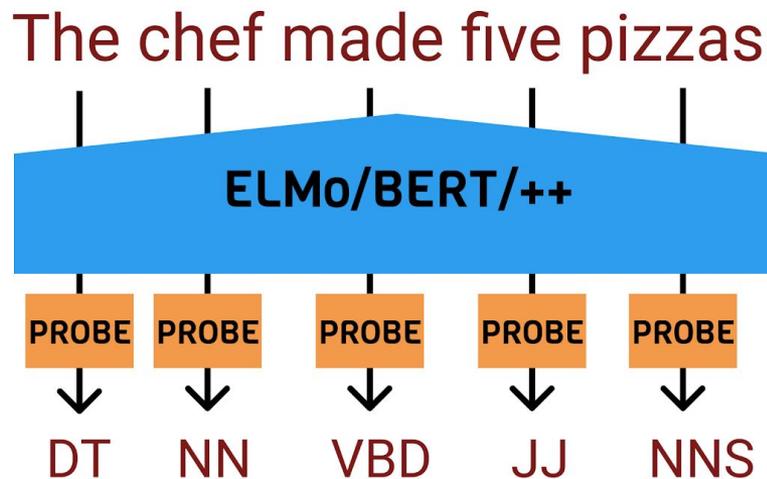


Probing: supervised analysis of representations

Does my network make task (e.g., part-of-speech) labels accessible?

Choose a function family to decode the task. (e.g., linear)

Train a function representations --> task



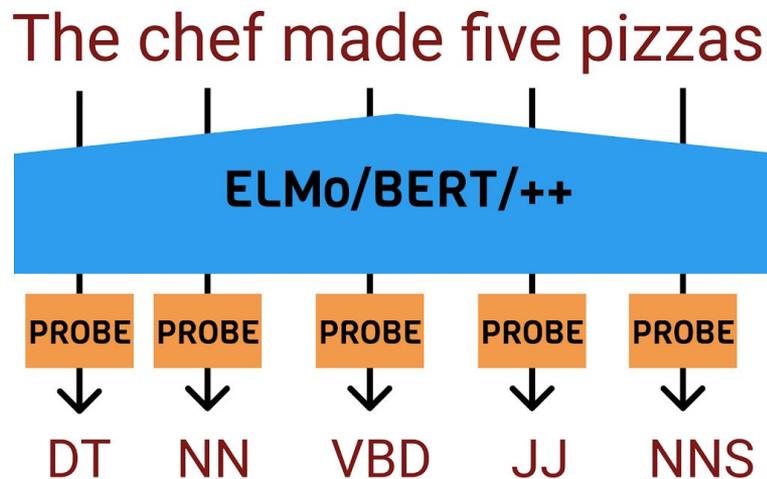
Probing: supervised analysis of representations

Does my network make task (e.g., part-of-speech) labels accessible?

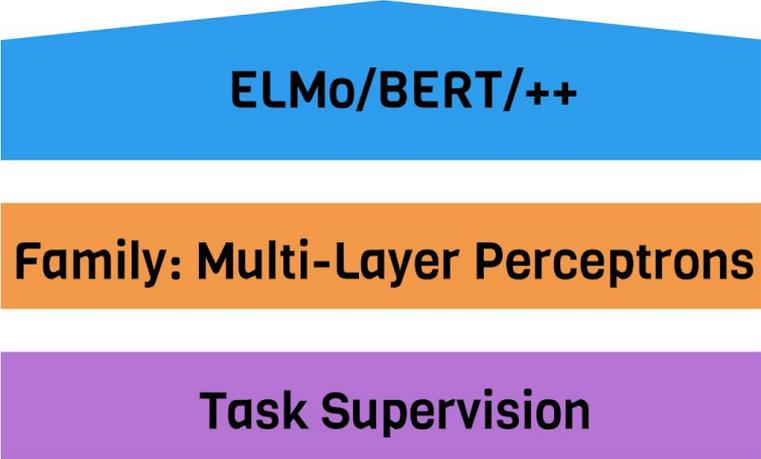
Choose a function family to decode the task. (e.g., linear)

Train a function representations --> task

Interpret accuracy on held-out data



The probe confounder problem

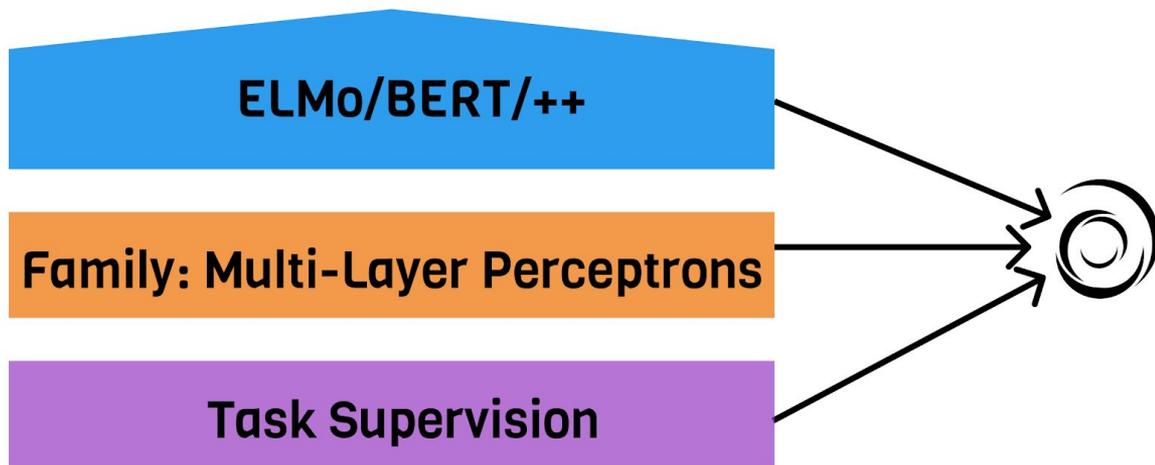


ELMo/BERT/++

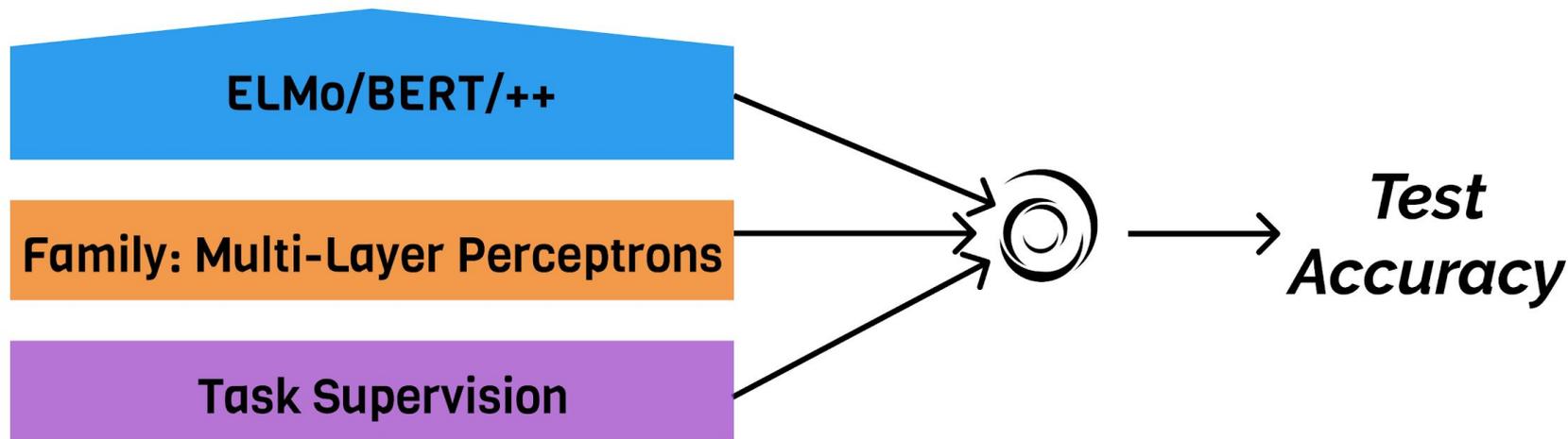
Family: Multi-Layer Perceptrons

Task Supervision

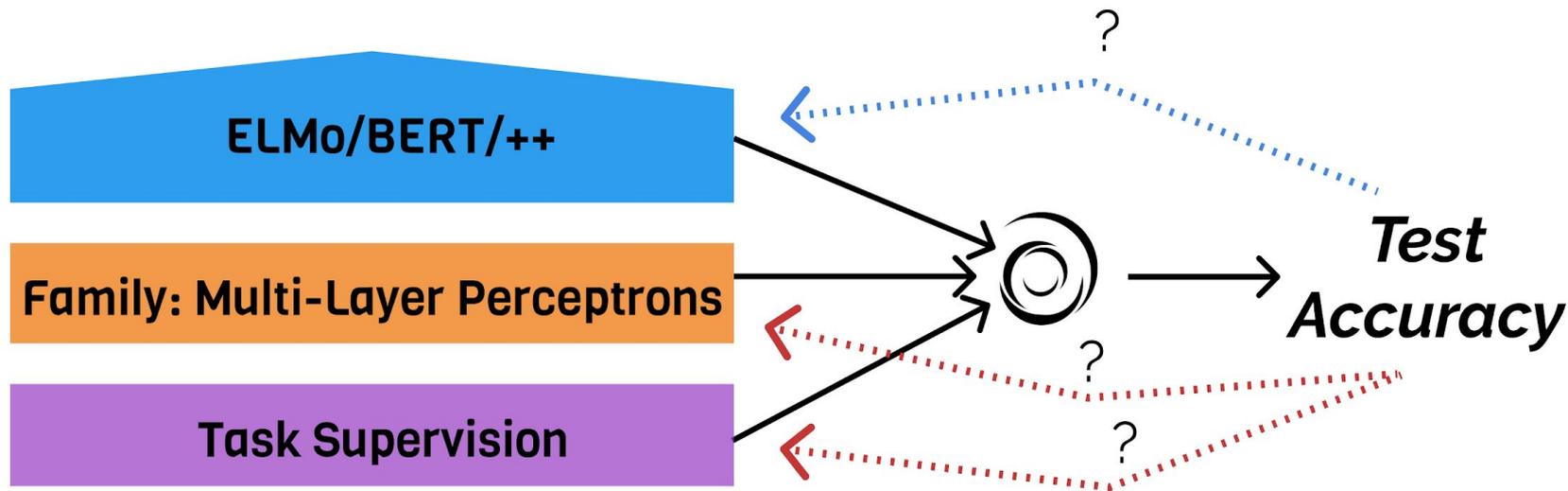
The probe confounder problem



The probe confounder problem



The probe confounder problem



Should we give credit to the **representation**?
(and/or) the **probe** and the **task supervision**?

This work:

This work:

1. Does high probe test accuracy mean the representation learned the task?

This work:

1. Does high probe test accuracy mean the representation learned the task?
2. How does the design of probes affect probing results?

This work:

1. Does high probe test accuracy mean the representation learned the task?
2. How does the design of probes affect probing results?
3. Can the probe confounder problem affect probing conclusions in practice?

Question 1

Does high probe test accuracy mean the representation learned a task?

No. Our *control tasks* are learned by probes but not encoded by representations.

Question 1

Does high probe test accuracy mean the representation learned a task?

No. Our *control tasks* are learned by probes but not encoded by representations.

Probing: Does ELMo learn part-of-speech?

Probing: Does ELMo learn part-of-speech?

Train and test probes on ELMo representations on the Penn Treebank

Probing: Does ELMo learn part-of-speech?

Train and test probes on ELMo representations on the Penn Treebank

Representation

ELMo, layer 1

Probing: Does ELMo learn part-of-speech?

Train and test probes on ELMo representations on the Penn Treebank

Representation

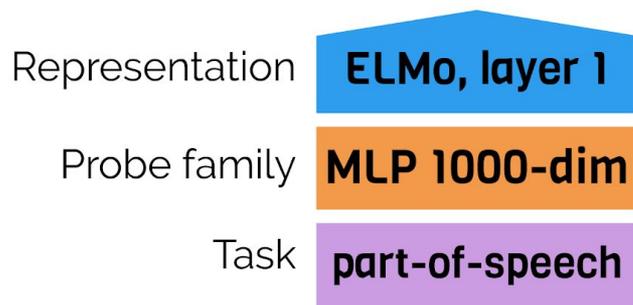
ELMo, layer 1

Probe family

MLP 1000-dim

Probing: Does ELMo learn part-of-speech?

Train and test probes on ELMo representations on the Penn Treebank



Probing: Does ELMo learn part-of-speech?

Train and test probes on ELMo representations on the Penn Treebank

Representation	ELMo, layer 1
Probe family	MLP 1000-dim
Task	part-of-speech
Test Accuracy	97.3

Probing: Does ELMo learn part-of-speech?

Train and test probes on ELMo representations on the Penn Treebank

Representation	ELMo, layer 1
Probe family	MLP 1000-dim
Task	part-of-speech
Test Accuracy	97.3

Probe achieves high accuracy!

Does the accuracy faithfully reflect the extent to which ELMo has learned part-of-speech tagging?

Defining **control tasks** for linguistic tasks

Defining **control tasks** for linguistic tasks

1. Look at task output space. e.g., 45 parts-of-speech.

Defining **control tasks** for linguistic tasks

1. Look at task output space. e.g., 45 parts-of-speech.
2. **Randomly partition** vocabulary into 45 categories

Defining **control tasks** for linguistic tasks

1. Look at task output space. e.g., 45 parts-of-speech.
2. **Randomly partition** vocabulary into 45 categories

Category1

house,
eat,

...

Defining **control tasks** for linguistic tasks

1. Look at task output space. e.g., 45 parts-of-speech.
2. **Randomly partition** vocabulary into 45 categories

Category1

house,
eat,

...

Category2

pizza,
the,

...

Defining **control tasks** for linguistic tasks

1. Look at task output space. e.g., 45 parts-of-speech.
2. **Randomly partition** vocabulary into 45 categories



Defining **control tasks** for linguistic tasks

1. Look at task output space. e.g., 45 parts-of-speech.
2. **Randomly partition** vocabulary into 45 categories



3. **Deterministically label** sentences in a corpus by looking up category for each word

Defining **control tasks** for linguistic tasks

1. Look at task output space. e.g., 45 parts-of-speech.
2. **Randomly partition** vocabulary into 45 categories

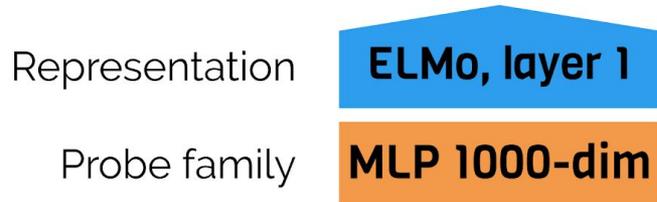


3. **Deterministically label** sentences in a corpus by looking up category for each word

the house is quiet as the people eat pizza

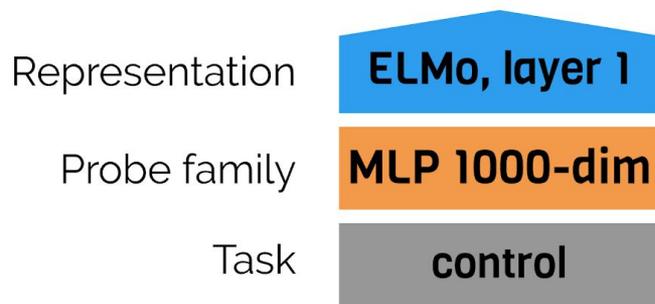
Do probes learn control tasks?

Train and test probes on ELMo representations on the Penn Treebank



Do probes learn control tasks?

Train and test probes on ELMo representations on the Penn Treebank



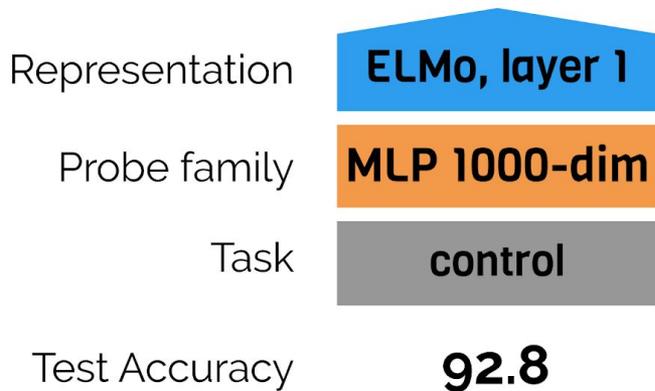
Do probes learn control tasks?

Train and test probes on ELMo representations on the Penn Treebank

Representation	ELMo, layer 1
Probe family	MLP 1000-dim
Task	control
Test Accuracy	92.8

Do probes learn control tasks?

Train and test probes on ELMo representations on the Penn Treebank



MLP probe: high accuracy on control tasks; does not reflect representation!

Do probes learn control tasks?

Train and test probes on ELMo representations on the Penn Treebank

Representation	ELMo, layer 1	ELMo, layer 1
Probe family	MLP 1000-dim	Linear
Task	control	control
Test Accuracy	92.8	

MLP probe: high accuracy on control tasks; does not reflect representation!

Do probes learn control tasks?

Train and test probes on ELMo representations on the Penn Treebank

Representation	ELMo, layer 1	ELMo, layer 1
Probe family	MLP 1000-dim	Linear
Task	control	control
Test Accuracy	92.8	71.2

MLP probe: high accuracy on control tasks; does not reflect representation!

Do probes learn control tasks?

Train and test probes on ELMo representations on the Penn Treebank

Representation	ELMo, layer 1	ELMo, layer 1
Probe family	MLP 1000-dim	Linear
Task	control	control
Test Accuracy	92.8	71.2

MLP probe: high accuracy on control tasks; does not reflect representation!

Linear probe: lower accuracy on control tasks

Selectivity for interpreting probing results

Idea: get a rough measure of how linguistic task accuracy may derive from probe expressivity and supervision.

Selectivity for interpreting probing results

Idea: get a rough measure of how linguistic task accuracy may derive from probe expressivity and supervision.

We define ***selectivity*** as a probe's accuracy on the linguistic task minus its accuracy on the control task

Probing part-of-speech vs control task

Can control tasks and selectivity help put probing accuracies in context?

Representation	ELMo, layer 1	ELMo, layer 1
Probe family	MLP 1000-dim	MLP 1000-dim
Task	control	part-of-speech
Test Accuracy	92.8	97.3

Probing part-of-speech vs control task

Can control tasks and selectivity help put probing accuracies in context?

Representation	ELMo, layer 1	ELMo, layer 1
Probe family	MLP 1000-dim	MLP 1000-dim
Task	control	part-of-speech
Test Accuracy	92.8	97.3
Selectivity		4.5

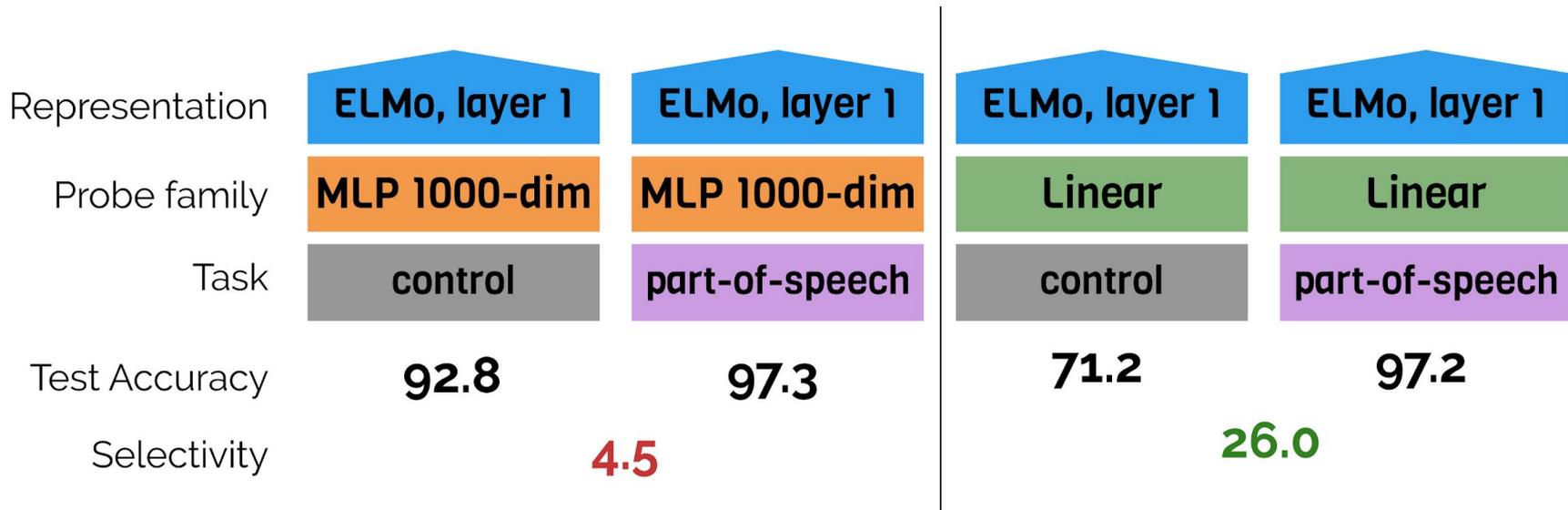
Probing part-of-speech vs control task

Can control tasks and selectivity help put probing accuracies in context?

Representation	ELMo, layer 1	ELMo, layer 1	ELMo, layer 1
Probe family	MLP 1000-dim	MLP 1000-dim	Linear
Task	control	part-of-speech	control
Test Accuracy	92.8	97.3	71.2
Selectivity		4.5	

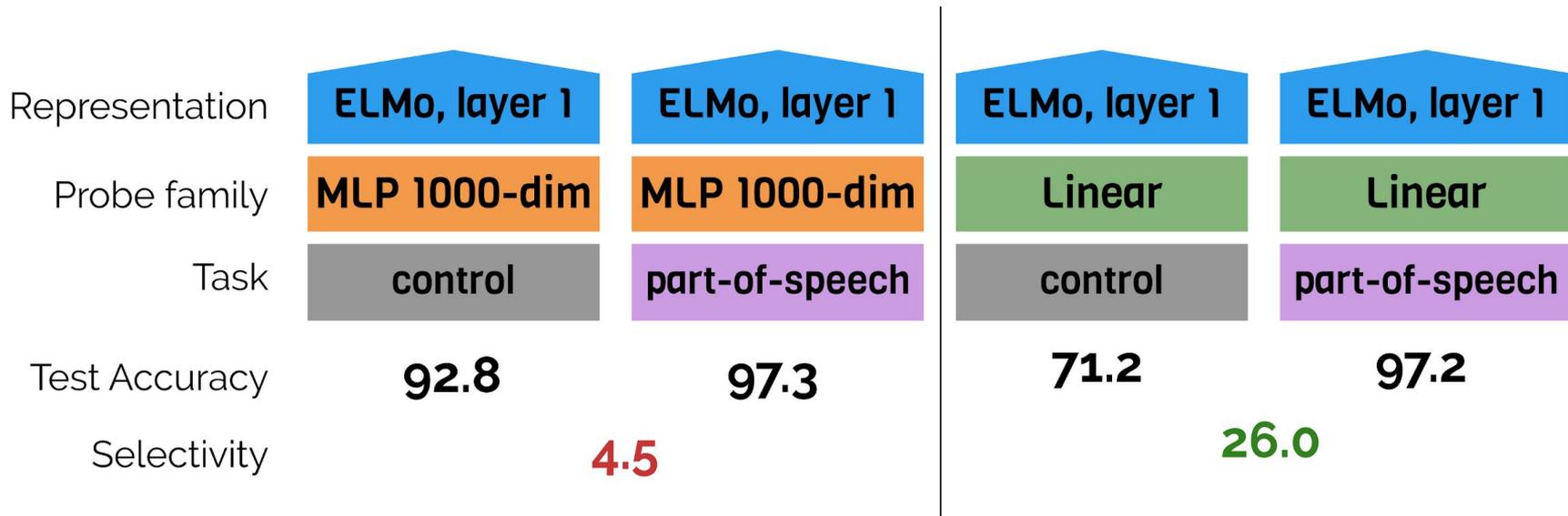
Probing part-of-speech vs control task

Can control tasks and selectivity help put probing accuracies in context?



Probing part-of-speech vs control task

Can control tasks and selectivity help put probing accuracies in context?



Probes with similar linguistic task accuracy may have very different selectivity

Question 2

How does the design of probes affect probing results?

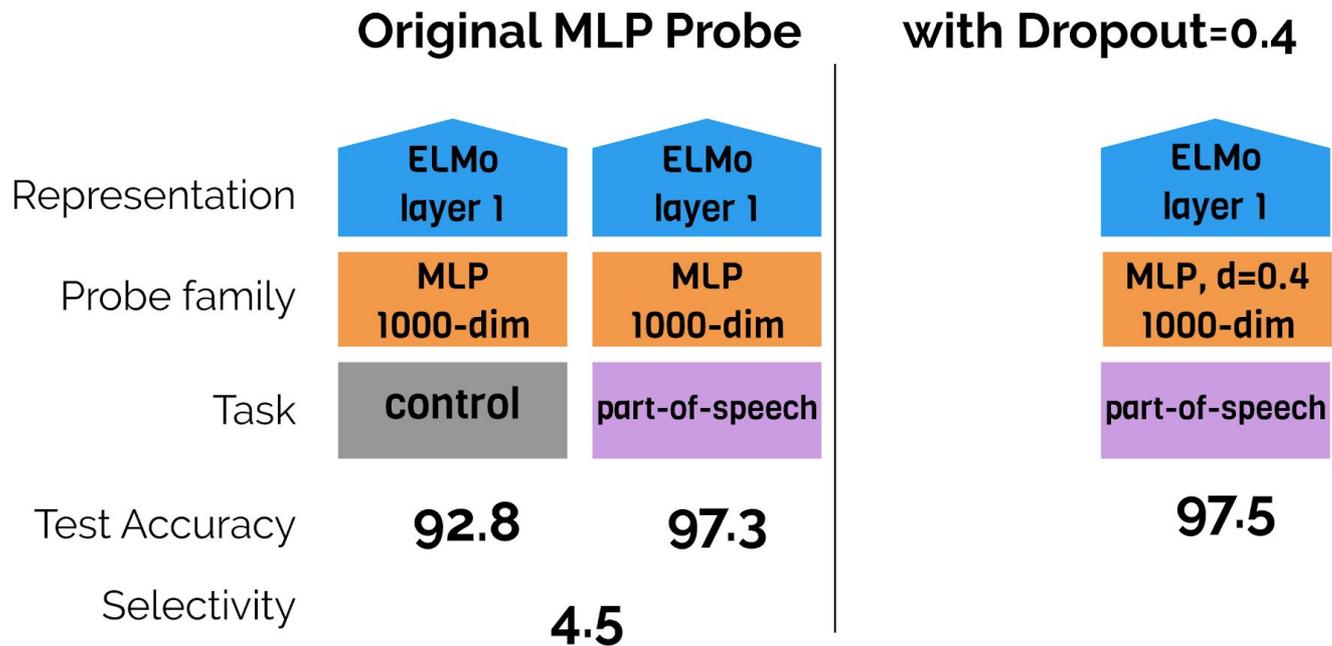
**Designing for good linguistic task generalization
does not necessarily lead to selective probes**

Designing probes with control tasks

Original MLP Probe

Representation	ELMo layer 1	ELMo layer 1
Probe family	MLP 1000-dim	MLP 1000-dim
Task	control	part-of-speech
Test Accuracy	92.8	97.3
Selectivity		4.5

Designing probes with control tasks



Designing probes with control tasks

	Original MLP Probe		with Dropout=0.4	
Representation	ELMo layer 1	ELMo layer 1	ELMo layer 1	ELMo layer 1
Probe family	MLP 1000-dim	MLP 1000-dim	MLP, d=0.4 1000-dim	MLP, d=0.4 1000-dim
Task	control	part-of-speech	control	part-of-speech
Test Accuracy	92.8	97.3	93.4	97.5
Selectivity		4.5		4.1

Designing probes with control tasks

	Original MLP Probe		with Dropout=0.4		with tiny hidden state (no dropout)	
Representation	ELMo layer 1	ELMo layer 1	ELMo layer 1	ELMo layer 1	ELMo layer 1	ELMo layer 1
Probe family	MLP 1000-dim	MLP 1000-dim	MLP, d=0.4 1000-dim	MLP, d=0.4 1000-dim	MLP 10-dim	MLP 10-dim
Task	control	part-of-speech	control	part-of-speech	control	part-of-speech
Test Accuracy	92.8	97.3	93.4	97.5	80.6	97.2
Selectivity		4.5		4.1		16.6

Designing probes with control tasks

	Original MLP Probe		with Dropout=0.4		with tiny hidden state (no dropout)	
Representation	ELMo layer 1	ELMo layer 1	ELMo layer 1	ELMo layer 1	ELMo layer 1	ELMo layer 1
Probe family	MLP 1000-dim	MLP 1000-dim	MLP, d=0.4 1000-dim	MLP, d=0.4 1000-dim	MLP 10-dim	MLP 10-dim
Task	control	part-of-speech	control	part-of-speech	control	part-of-speech
Test Accuracy	92.8	97.3	93.4	97.5	80.6	97.2
Selectivity		4.5		4.1		16.6

Simply regularizing — to minimize generalization gap — doesn't necessarily lead to selectivity!

Question 3

Can the probe confounder problem affect probing conclusions in practice?

Yes — probes may be picking up on spurious signals

Re-examining probes on ELMo's layers

Is ELMo1 better at part-of-speech than ELMo2?

Representation

ELMo, layer 1

Probe family

Linear

Task

part-of-speech

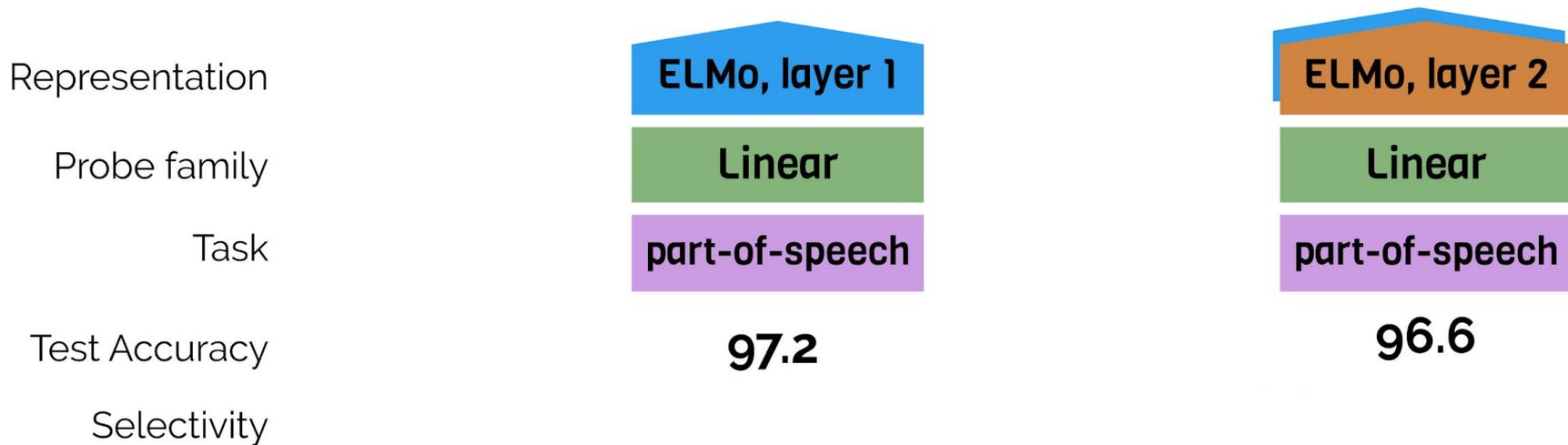
Test Accuracy

97.2

Selectivity

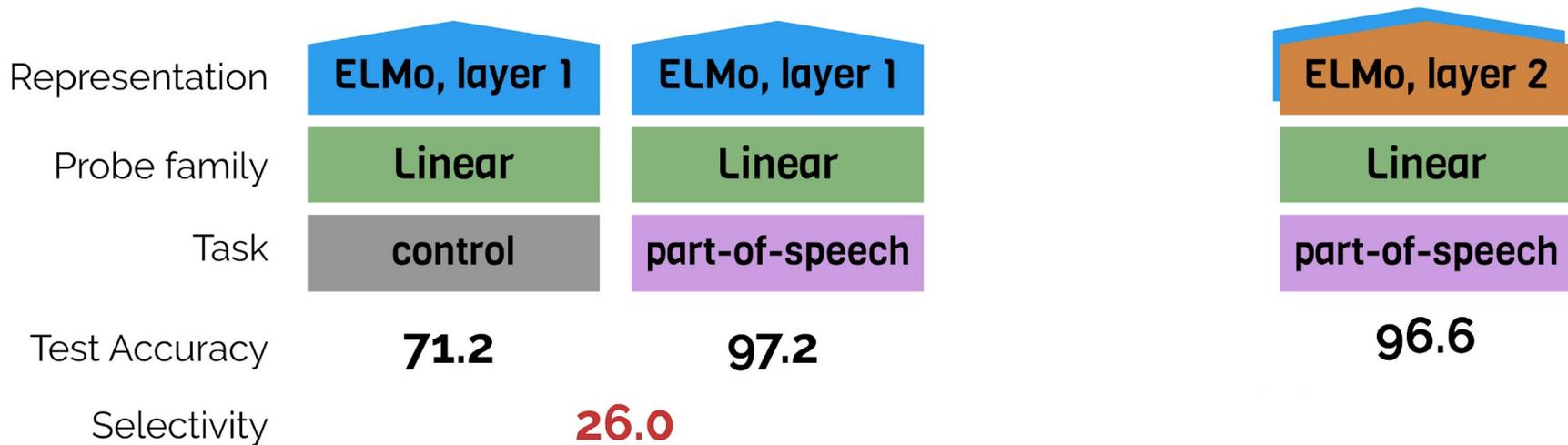
Re-examining probes on ELMo's layers

Is ELMo1 better at part-of-speech than ELMo2?



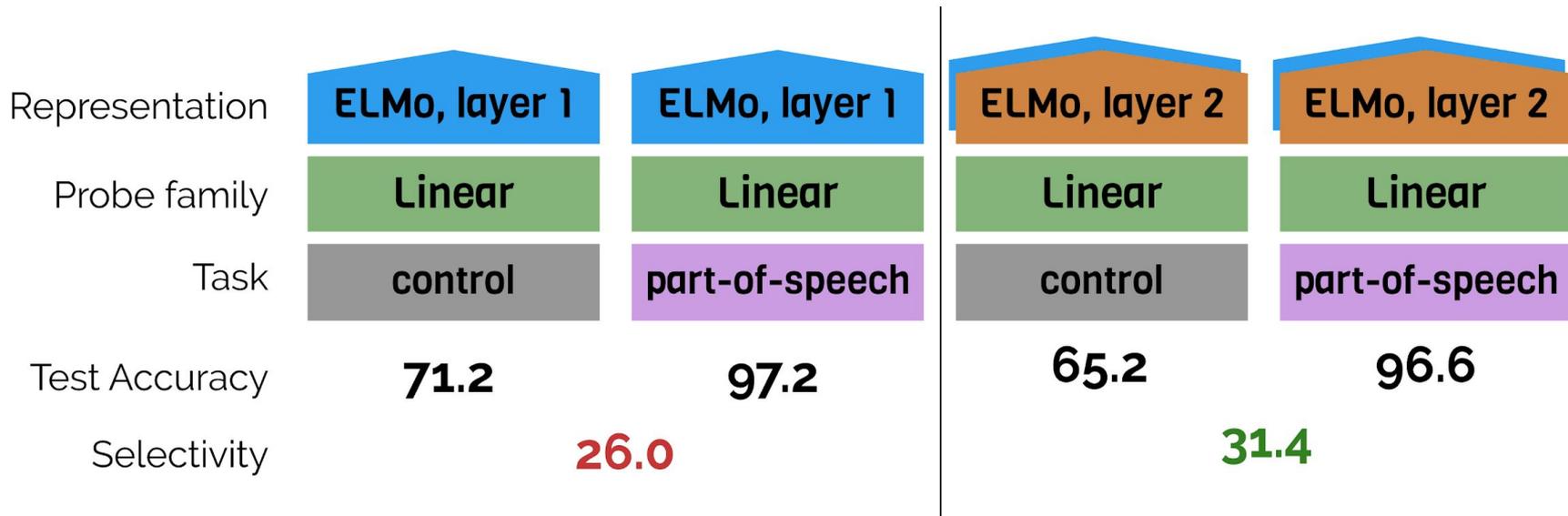
Re-examining probes on ELMo's layers

Is ELMo1 better at part-of-speech than ELMo2?



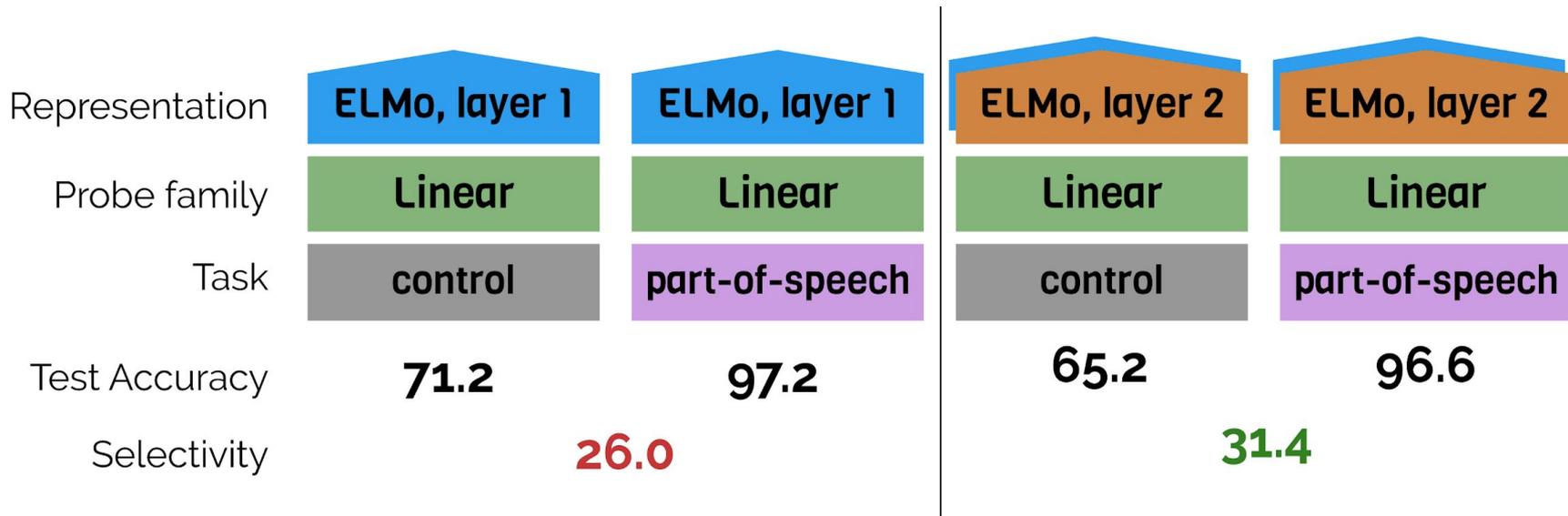
Re-examining probes on ELMo's layers

Is ELMo1 better at part-of-speech than ELMo2?



Re-examining probes on ELMo's layers

Is ELMo1 better at part-of-speech than ELMo2?



ELMo1 part-of-speech gains over ELMo2 may be explained by easier access to a **spurious signal: word identity**

Limitations

Limitations

Our control tasks only use **word identity**; there are many possible spurious signals in probing

Limitations

Our control tasks only use **word identity**; there are many possible spurious signals in probing

Selectivity **builds intuition** but does not permit fine-grained claims, like “my model got *this* selectivity, so it learned the task.”

Thanks!

