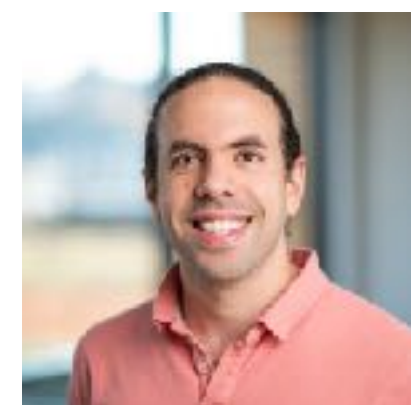


Transparency Helps Reveal When Language Models Learn Meaning

TACL 2023

Zhaofeng Wu, Will Merrill, Hao Peng, Iz Beltagy, and Noah Smith



Δ⁺ <C^oσ> ἰ⁺ἰ⁺Δ⁺ Δ⁺> ἰ⁺ἰ⁺Δ⁺ Ἀ⁺ἰ⁺ἰ⁺ Δ⁺. Δ⁺σ⁺ ἰ⁺ἰ⁺Δ⁺
Ἀ⁺ἰ⁺ ἰ⁺ἰ⁺ <ἰ⁺ἰ⁺ἰ⁺ ἰ⁺ἰ⁺ Ἀ⁺ἰ⁺ Ἀ⁺ἰ⁺ἰ⁺. ἰ⁺ἰ⁺ Ἀ⁺ἰ⁺ ἰ⁺ἰ⁺ἰ⁺ Ἀ⁺ἰ⁺ἰ⁺
<ἰ⁺ἰ⁺ἰ⁺. Δ⁺ἰ⁺ἰ⁺ ἰ⁺ἰ⁺ἰ⁺

ᐃᐅ ᐊᑕᓐᓂᐅᑦ ᐱᑕᓴᐊᑦᐅᑦ ᐃᓂᓴᐅᓴᐅᑦ ᐱᓴᐱᑦᐅᑦ ᐅᓄᓄᓄᓴᐅᑦ ᐅᓄᓴᓂᓴᐅᑦ ᐱᑕᓴᐊᑦᐅᑦ ᐱᑦᐅᑦ ᐱᓴᐅᑦ ᐊᑕᓂᐊᑦᐅᑦ ᐊᑦᐅᑦ ᐱᑦᐅᑦ ᐊᓴᐅᓂᑦᑦ ᐱᓴᐅᑦᐅᑦ ᐊᑕᓂᐊᑦᐅᑦ ᐃᓴᐊᓴᐅᑦ ᑕᓴᐅᓴᑕᓴᐅᑦ



LMs can't learn meaning from form alone.



LMs can't learn meaning from form alone.

Can we say LMs understand language?


```
def f(n):  
    if n == 1 or n == 2:  
        return 1  
    return f(n - 1) + f(n - 2)
```



```
def f(n):  
    if n == 1 or n == 2:  
        return 1  
    return f(n - 1) + f(n - 2)
```



LMs can't learn execution.

```
def f(n):  
    if n == 1 or n == 2:  
        return 1  
    return f(n - 1) + f(n - 2)
```



LMs can't learn execution.

There are assertions:
assert f(6) == 8




```
def f(n):
    if n == 1 or n == 2:
        return 1
    return f(n - 1) + f(n - 2)
```

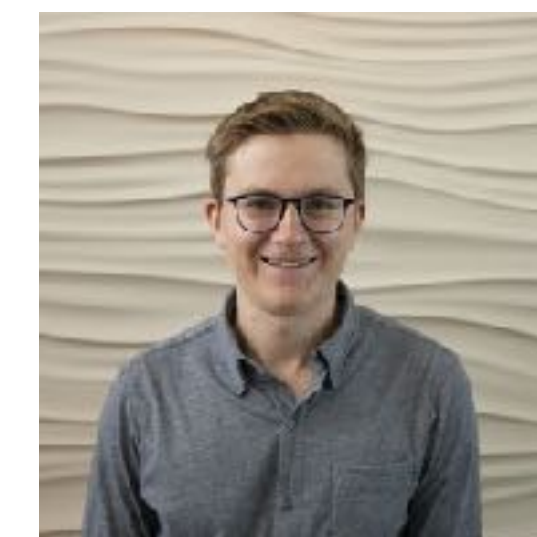



LMs can't learn execution.

There are assertions:
assert f(6) == 8



Assertions enable meaning learnability in some languages.



Google Scholar		
 et al. The academic superstar everybody wants to be co-author with.		
Cited by	VIEW ALL	
	All	Since 2017
Citations	3700948	955667
h-index	333	250
i10-index	333	333

```
def f(n):
    if n == 1 or n == 2:
        return 1
    return f(n - 1) + f(n - 2)
```

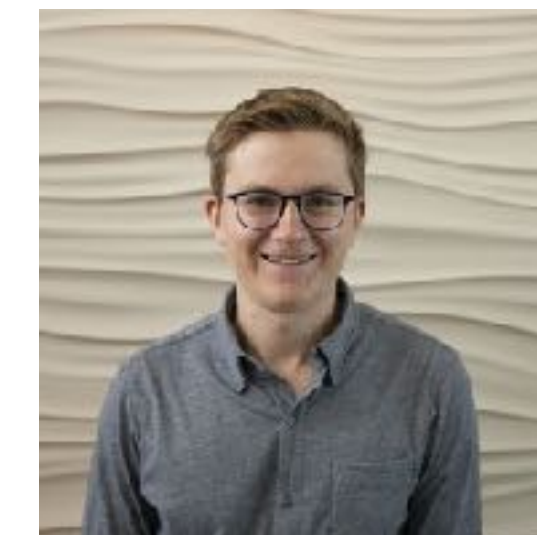


LMs can't learn execution.

There are assertions:
assert f(6) == 8



Assertions enable meaning learnability in some languages.



Google Scholar		
et al. The academic superstar everybody wants to be co-author with.		
Cited by	VIEW ALL	
	All	Since 2017
Citations	3700948	955667
h-index	333	250
i10-index	333	333

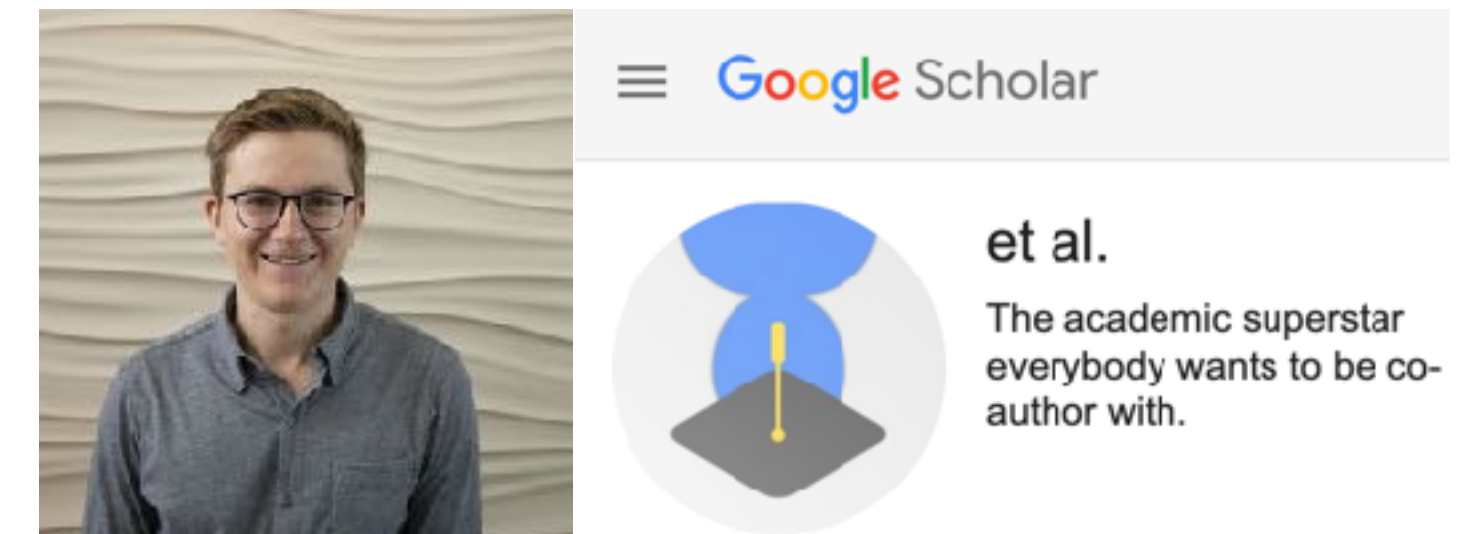


LMs can't learn execution.

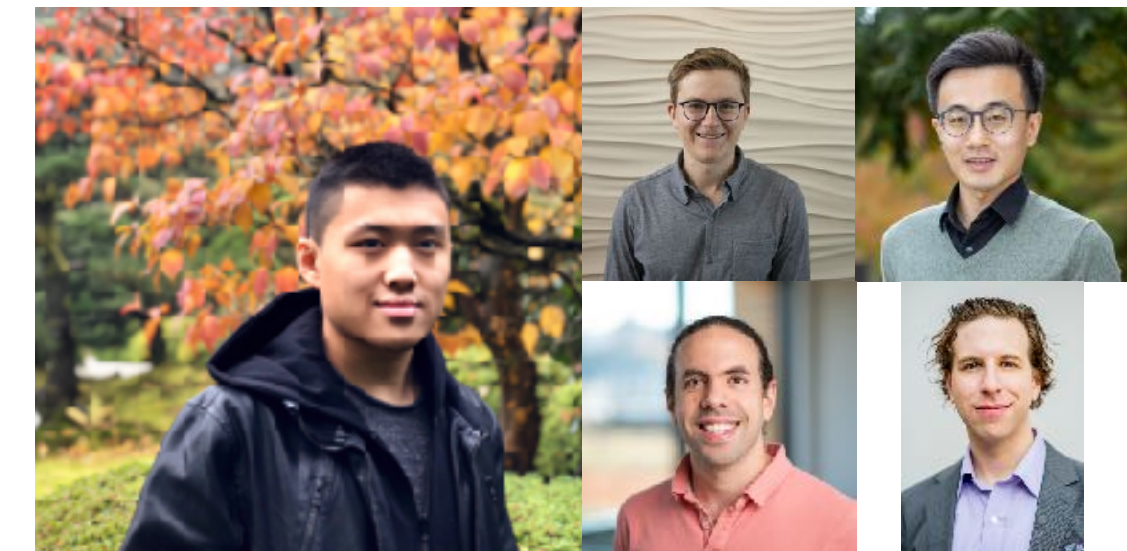
There are assertions:
`assert f(6) == 8`



Assertions enable meaning learnability in some languages.



LMs learn the meaning of some languages with assertions.



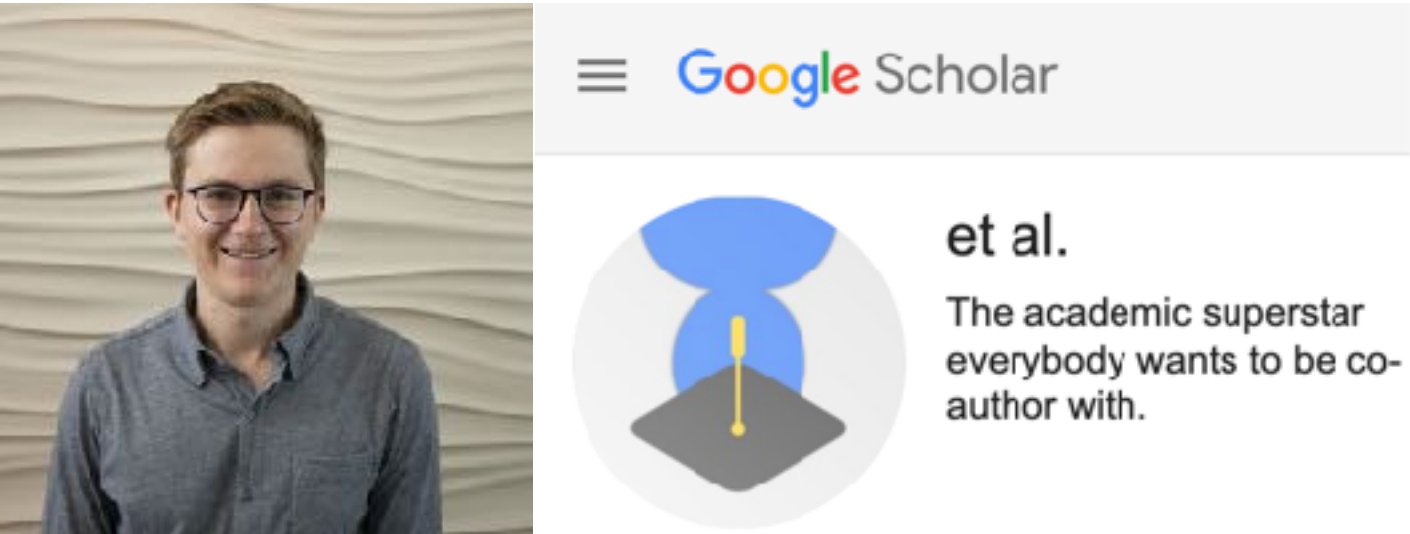


LMs can't learn execution.

There are assertions:
`assert f(6) == 8`

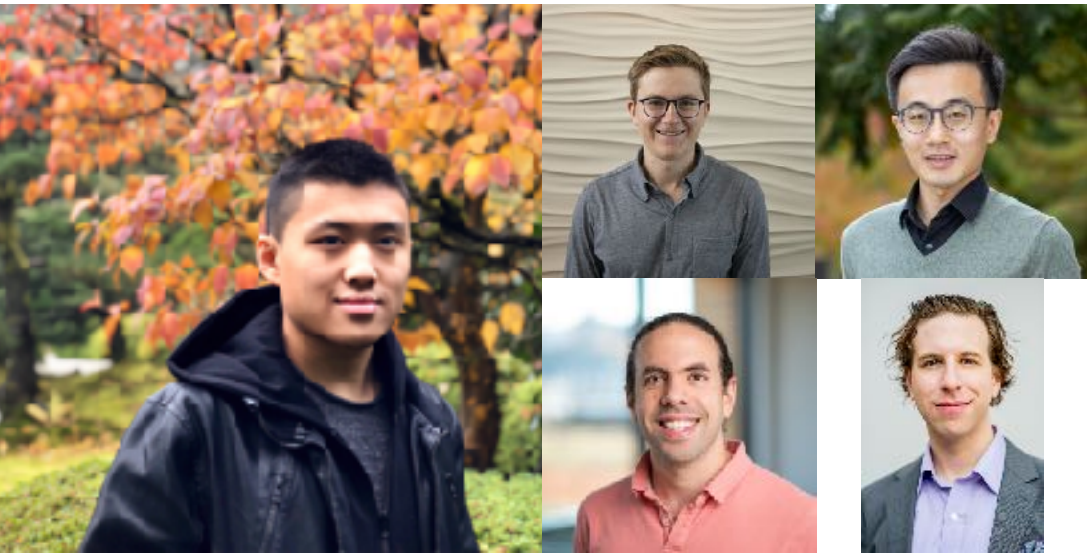


Assertions enable meaning learnability in some languages.



LMs learn the meaning of some languages with assertions.

But not natural language.



Can LMs Learn From Assertions?

Setup

Setup

```
((¬T)∧(¬(Tv(¬F))))=(Tv(¬(¬((¬T)∨(¬(¬F)))))  
¬(¬(¬((F∧((F∧F)∧F))∧F)∧(¬T))))=(T∧T)∧((¬F)∨(¬F))  
(((¬(¬(¬(¬(¬T))))∨T)∨T)∧(¬(¬T)))=(¬F)∨(¬(T∧(TvT)))  
(T∧(FvF))∨(Tv(F∧T))=¬((¬T)∧(¬(¬(¬(¬F))vF))v(T∧T)))  
(((¬(¬F))∧(¬F))∧((¬F)vF)∧F)=(F∧(¬(¬(Fv(¬(Fv(¬T)))∧T))))  
(Tv(¬(T∧(Tv(¬(Fv(¬F)))))∨T)=¬((¬(Tv(¬(¬(¬(T∧F)))))∧F))  
¬((¬(¬F))v((¬F)v(Tv(¬(TvT)))))=((F∧(¬T))∧((¬F)∧F))∧F  
(F∧(F∧(¬(¬(TvT))∧(¬T))))=¬(¬((¬T)vF)v(Fv(¬(¬F))))  
(F∧(F∧(¬((FvF)v(¬(¬T)))))=¬(((¬(T∧T))v(¬F))v(¬T))∧(¬F))
```

Pretraining

Setup

$((\neg T) \wedge (\neg(T \vee (\neg F)))) = (T \vee (\neg(\neg((\neg T) \vee (\neg(\neg F)))))$
 $(\neg(\neg(\neg((F \wedge ((F \wedge F) \wedge F)) \wedge F) \wedge (\neg T)))) = ((T \wedge T) \wedge ((\neg F) \vee (\neg F)))$
 $((\neg(\neg(\neg(\neg(\neg T)))) \vee T) \vee T) \wedge (\neg(\neg T)) = ((\neg F) \vee (\neg(T \wedge (T \vee T))))$
 $((T \wedge (F \vee F)) \vee (T \vee (F \wedge T))) = (\neg((\neg T) \wedge (\neg(\neg(\neg(\neg F)) \vee F)) \vee (T \wedge T)))$
 $((\neg(\neg F)) \wedge (\neg F)) \wedge ((\neg F) \vee F) \wedge F) = (F \wedge (\neg(\neg(F \vee (\neg(F \vee (\neg T))) \wedge T))))$
 $((T \vee (\neg(T \wedge (T \vee (\neg(F \vee (\neg F))))) \vee T) = (\neg((\neg(T \vee (\neg(\neg(\neg(T \wedge F))))) \wedge F))$
 $(\neg((\neg(\neg F)) \vee ((\neg F) \vee (T \vee (\neg(T \vee T))))) = (((F \wedge (\neg T)) \wedge ((\neg F) \wedge F)) \wedge F)$
 $(F \wedge (F \wedge (\neg((\neg(T \vee T)) \wedge (\neg T)))) = (\neg(\neg((\neg T) \vee F) \vee (F \vee (\neg(\neg F)))))$
 $(F \wedge (F \wedge (\neg((F \vee F) \vee (\neg(\neg T))))) = (\neg(((\neg(T \wedge T)) \vee (\neg F)) \vee (\neg T)) \wedge (\neg F))$



RoBERTa-like MLM
GPT-2-like ALM

Pretraining

Setup

```
((¬T)∧(¬(Tv(¬F))))=(Tv(¬(¬((¬T)∨(¬(¬F))))))
(¬(¬(¬((F∧((F∧F)∧F))∧F)∧(¬T))))=(T∧T)∧((¬F)∨(¬F))
(((¬((¬(¬(¬(¬T))))∨T))∨T)∧(¬(¬T)))=(¬F)∨(¬(T∧(TvT))))
((T∧(FvF))∨(Tv(F∧T)))=(¬((¬T)∧(¬((¬(¬(¬F))vF))v(T∧T))))
(((¬(¬F))∧(¬F))∧((¬F)∨F)∧F)=(F∧(¬(¬(Fv(¬(Fv(¬T)))∧T))))
((Tv(¬(T∧(Tv(¬(Fv(¬F))))))∨T)=(¬((¬(Tv(¬(¬(¬(T∧F))))))∧F))
(¬((¬(¬F))v((¬F)∨(Tv(¬(TvT))))))=((F∧(¬T))∧((¬F)∧F)∧F)
(F∧(F∧(¬((¬(TvT))∧(¬T))))=(¬(¬((¬T)∨F)∨(Fv(¬(¬F))))))
(F∧(F∧(¬((FvF)∨(¬(¬T))))))=(¬(((¬(T∧T))∨(¬F))∨(¬T))∧(¬F))
```

Pretraining



RoBERTa-like MLM
GPT-2-like ALM

Probing

Setup

$((\neg T) \wedge (\neg (T \vee (\neg F)))) = (T \vee (\neg (\neg ((\neg T) \vee (\neg (\neg F)))))$
 $(\neg (\neg (\neg ((F \wedge ((F \wedge F) \wedge F)) \wedge F) \wedge (\neg T)))) = ((T \wedge T) \wedge ((\neg F) \vee (\neg F)))$
 $((\neg (\neg (\neg (\neg (\neg T)))) \vee T) \vee T) \wedge (\neg (\neg T)) = ((\neg F) \vee (\neg (T \wedge (T \vee T))))$
 $((T \wedge (F \vee F)) \vee (T \vee (F \wedge T))) = (\neg ((\neg T) \wedge (\neg (\neg ((\neg (\neg F)) \vee F)) \vee (T \wedge T))))$
 $((\neg (\neg F)) \wedge (\neg F)) \wedge ((\neg F) \vee F) \wedge F) = (F \wedge (\neg (\neg (F \vee (\neg (F \vee (\neg T))) \wedge T))))$
 $((T \vee (\neg (T \wedge (T \vee (\neg (F \vee (\neg F))))) \vee T) = (\neg ((\neg (T \vee (\neg (\neg (\neg (T \wedge F))))) \wedge F)))$
 $(\neg ((\neg (\neg F)) \vee ((\neg F) \vee (T \vee (\neg (T \vee T))))) = (((F \wedge (\neg T)) \wedge ((\neg F) \wedge F)) \wedge F)$
 $(F \wedge (F \wedge (\neg ((\neg (T \vee T)) \wedge (\neg T)))) = (\neg (\neg ((\neg T) \vee F) \vee (F \vee (\neg (\neg F)))))$
 $(F \wedge (F \wedge (\neg ((F \vee F) \vee (\neg (\neg T))))) = (\neg (((\neg (T \wedge T)) \vee (\neg F)) \vee (\neg T)) \wedge (\neg F)))$

Pretraining



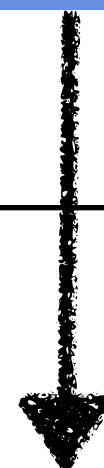
RoBERTa-like MLM
GPT-2-like ALM

Probing

Setup

$((\neg T) \wedge (\neg (T \vee (\neg F)))) = (T \vee (\neg (\neg ((\neg T) \vee (\neg (\neg F)))))$
 $(\neg (\neg (\neg ((F \wedge ((F \wedge F) \wedge F)) \wedge F) \wedge (\neg T)))) = ((T \wedge T) \wedge ((\neg F) \vee (\neg F)))$
 $((\neg (\neg (\neg (\neg (\neg T)))) \vee T) \vee T) \wedge (\neg (\neg T)) = ((\neg F) \vee (\neg (T \wedge (T \vee T))))$
 $((T \wedge (F \vee F)) \vee (T \vee (F \wedge T))) = (\neg ((\neg T) \wedge (\neg (\neg (\neg (\neg F)) \vee F)) \vee (T \wedge T))))$
 $((\neg (\neg F)) \wedge (\neg F)) \wedge ((\neg F) \vee F) \wedge F) = (F \wedge (\neg (\neg (F \vee (\neg (F \vee (\neg T))) \wedge T))))$
 $((T \vee (\neg (T \wedge (T \vee (\neg (F \vee (\neg F))))) \vee T) = (\neg ((\neg (T \vee (\neg (\neg (\neg (T \wedge F))))) \wedge F)))$
 $(\neg ((\neg (\neg F)) \vee ((\neg F) \vee (T \vee (\neg (T \vee T))))) = (((F \wedge (\neg T)) \wedge ((\neg F) \wedge F)) \wedge F)$
 $(F \wedge (F \wedge (\neg ((\neg (T \vee T)) \wedge (\neg T)))) = (\neg (\neg ((\neg T) \vee F) \vee (F \vee (\neg (\neg F)))))$
 $(F \wedge (F \wedge (\neg ((F \vee F) \vee (\neg (\neg T))))) = (\neg (((\neg (T \wedge T)) \vee (\neg F)) \vee (\neg T)) \wedge (\neg F)))$

Pretraining



$(T \vee (F \wedge T))$

$(F \wedge (\neg T))$

unseen



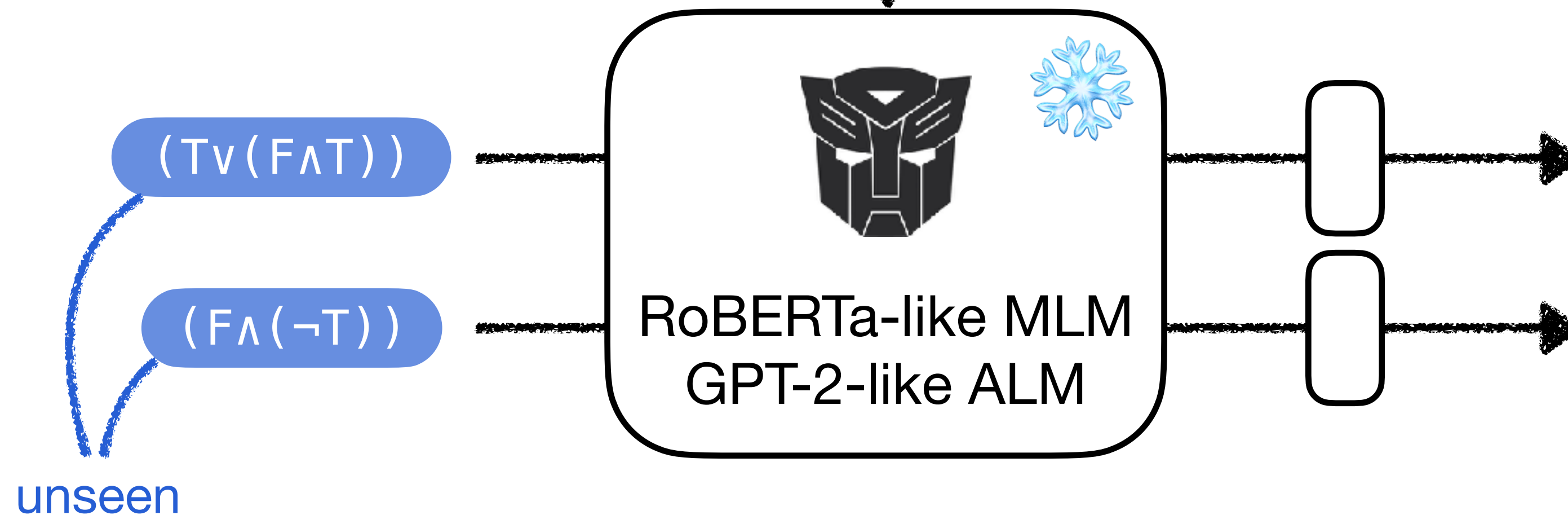
RoBERTa-like MLM
GPT-2-like ALM

Probing

Setup

```
((¬T)∧(¬(Tv(¬F))))=(Tv(¬(¬((¬T)∨(¬(¬F))))))
(¬(¬(¬((F∧((F∧F)∧F))∧F)∧(¬T))))=(T∧T)∧((¬F)∨(¬F))
(((¬((¬(¬(¬(¬T))))∨T))∨T)∧(¬(¬T)))=(¬F)∨(¬(T∧(TvT))))
((T∧(FvF))∨(Tv(F∧T)))=(¬((¬T)∧(¬((¬(¬(¬F))vF))v(T∧T))))
(((¬(¬F))∧(¬F))∧((¬F)∨F)∧F)=(F∧(¬(¬(Fv(¬(Fv(¬T)))∧T))))
((Tv(¬(T∧(Tv(¬(Fv(¬F))))))∨T)=(¬((¬(Tv(¬(¬(¬(T∧F))))))∧F))
(¬((¬(¬F))v((¬F)∨(Tv(¬(TvT))))))=((F∧(¬T))∧((¬F)∧F)∧F)
(F∧(F∧(¬((¬(TvT))∧(¬T))))=(¬(¬((¬T)∨F)∨(Fv(¬(¬F))))))
(F∧(F∧(¬((FvF)∨(¬(¬T))))))=(¬(((¬(T∧T))∨(¬F))∨(¬T))∧(¬F))
```

Pretraining



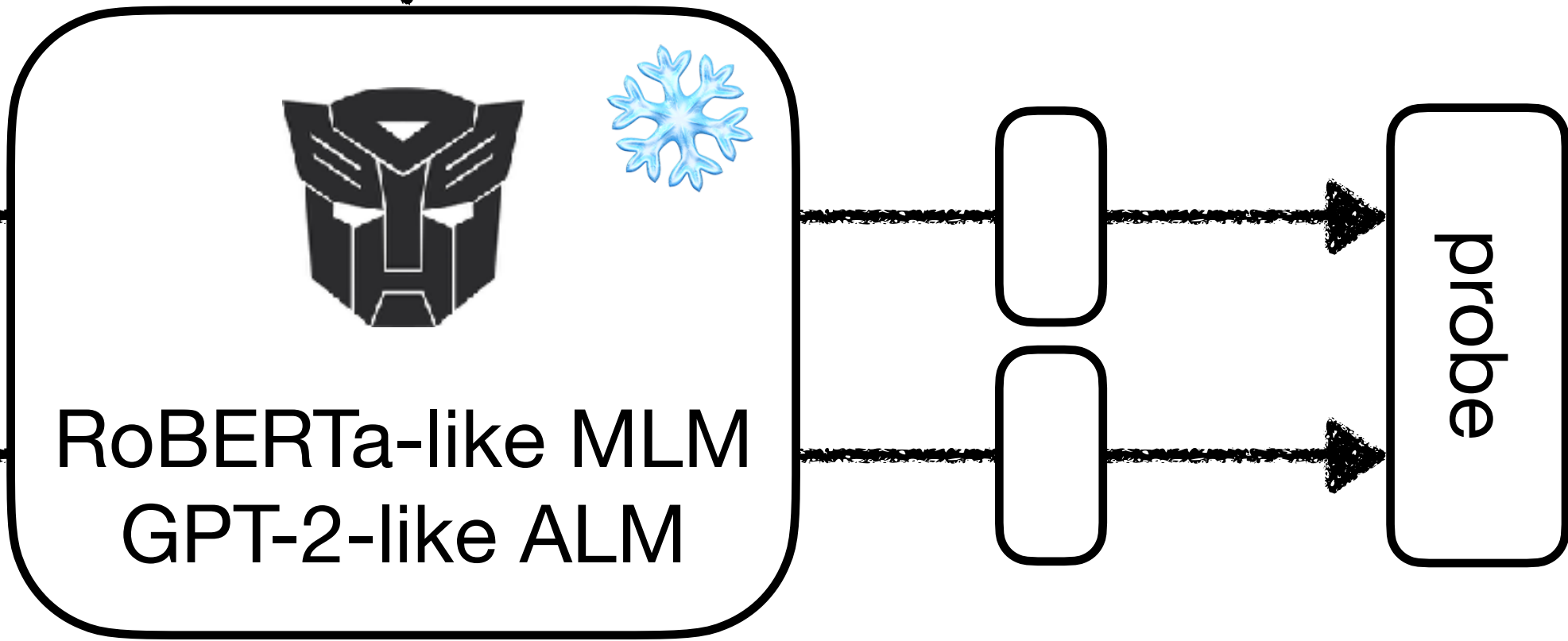
Setup

$((\neg T) \wedge (\neg(T \vee (\neg F)))) = (T \vee (\neg(\neg((\neg T) \vee (\neg(\neg F)))))$
 $(\neg(\neg(\neg((F \wedge ((F \wedge F) \wedge F)) \wedge F) \wedge (\neg T)))) = ((T \wedge T) \wedge ((\neg F) \vee (\neg F)))$
 $((\neg(\neg(\neg(\neg(\neg T)))) \vee T) \vee T) \wedge (\neg(\neg T)) = ((\neg F) \vee (\neg(T \wedge (T \vee T))))$
 $((T \wedge (F \vee F)) \vee (T \vee (F \wedge T))) = (\neg((\neg T) \wedge (\neg(\neg(\neg(\neg F)) \vee F)) \vee (T \wedge T)))$
 $((\neg(\neg F)) \wedge (\neg F)) \wedge ((\neg F) \vee F) \wedge F) = (F \wedge (\neg(\neg(F \vee (\neg(F \vee (\neg T))) \wedge T))))$
 $((T \vee (\neg(T \wedge (T \vee (\neg(F \vee (\neg F))))) \vee T) = (\neg((\neg(T \vee (\neg(\neg(\neg(T \wedge F))))) \wedge F))$
 $(\neg((\neg(\neg F)) \vee ((\neg F) \vee (T \vee (\neg(T \vee T))))) = (((F \wedge (\neg T)) \wedge ((\neg F) \wedge F)) \wedge F)$
 $(F \wedge (F \wedge (\neg((\neg(T \vee T)) \wedge (\neg T)))) = (\neg(\neg((\neg T) \vee F) \vee (F \vee (\neg(\neg F)))))$
 $(F \wedge (F \wedge (\neg((F \vee F) \vee (\neg(\neg T))))) = (\neg(((\neg(T \wedge T)) \vee (\neg F)) \vee (\neg T)) \wedge (\neg F))$

Pretraining

$(T \vee (F \wedge T))$
 $(F \wedge (\neg T))$

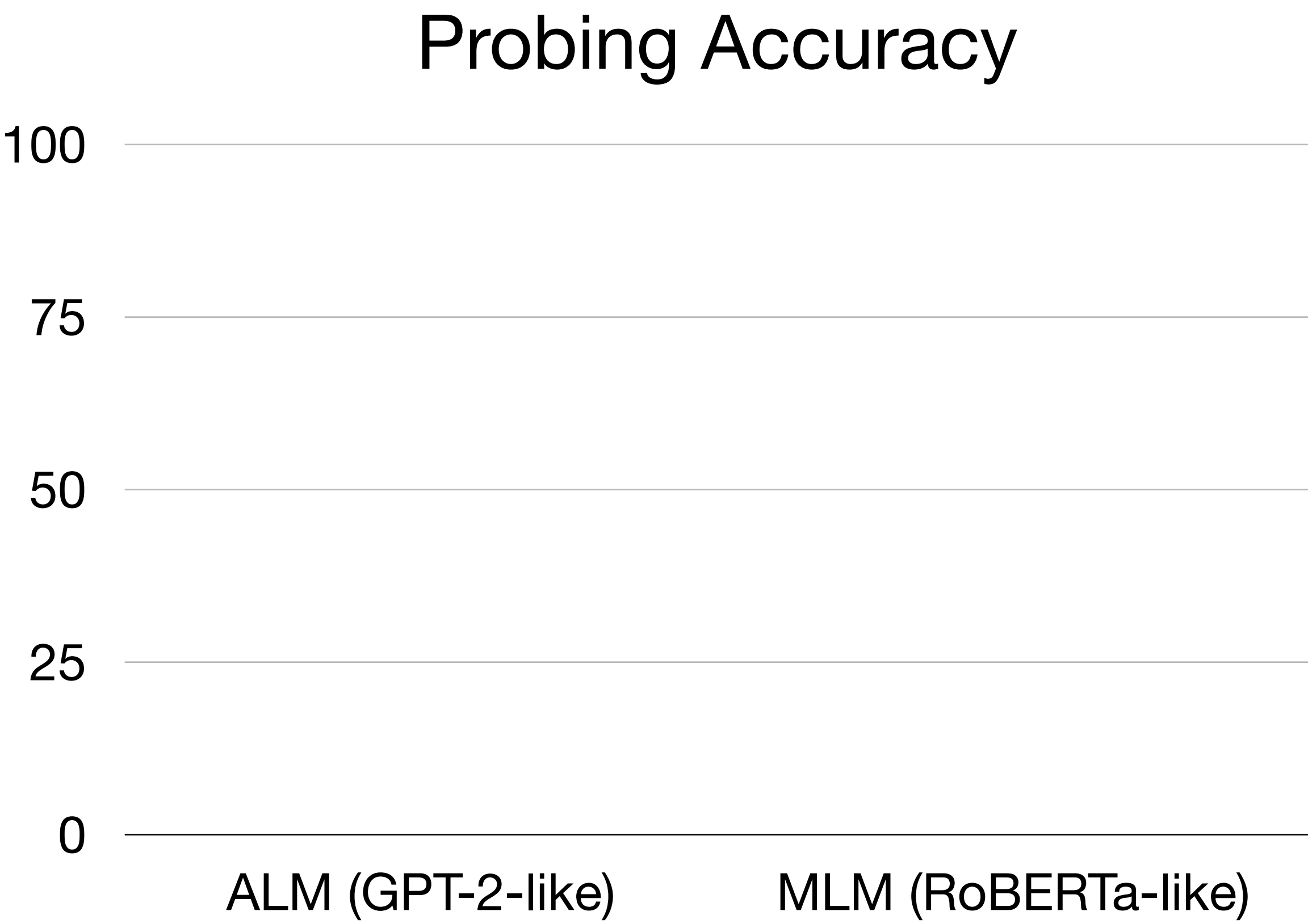
unseen



$\in \{ = , \neq \}$ Probing

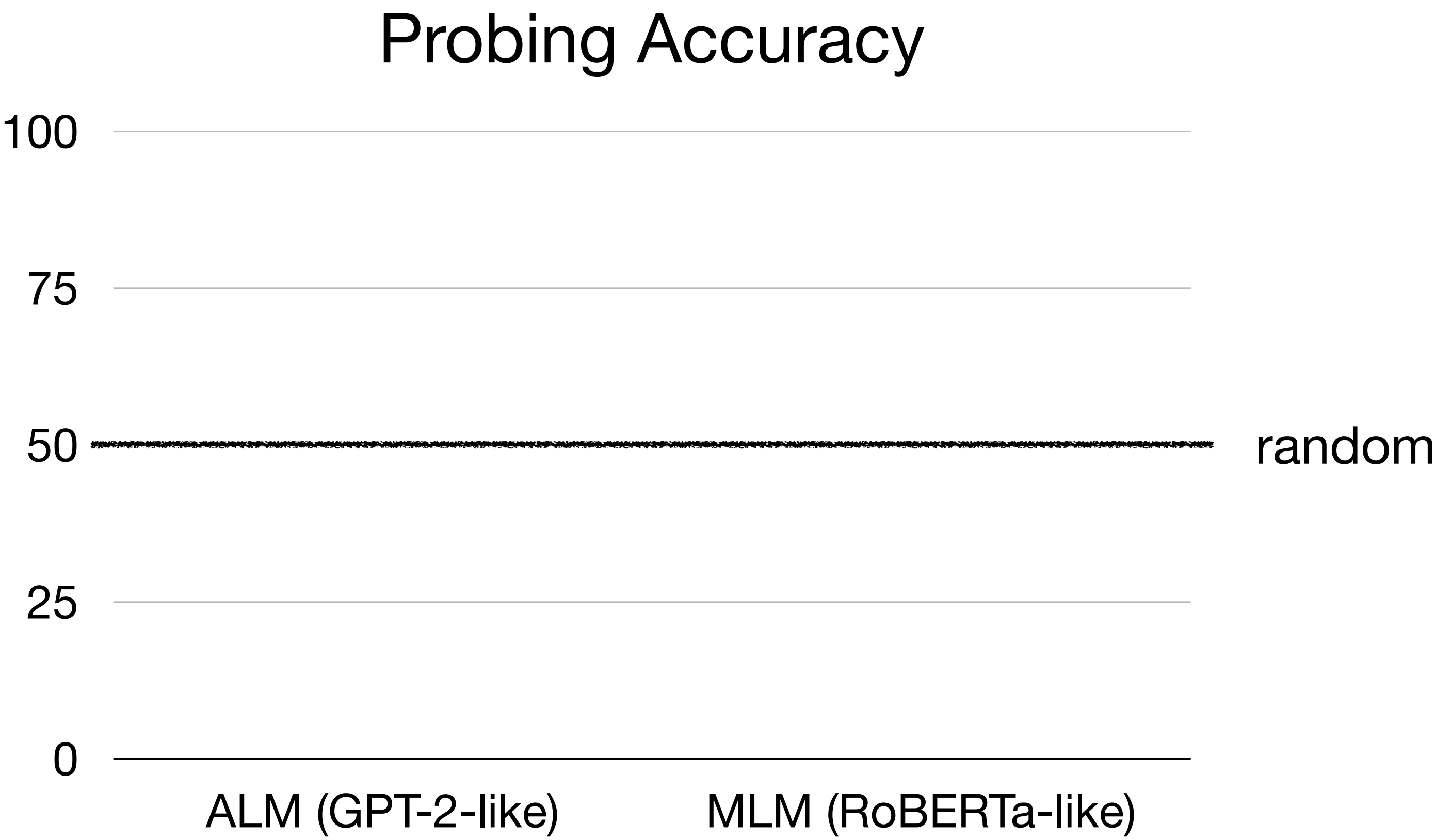
Results

Results



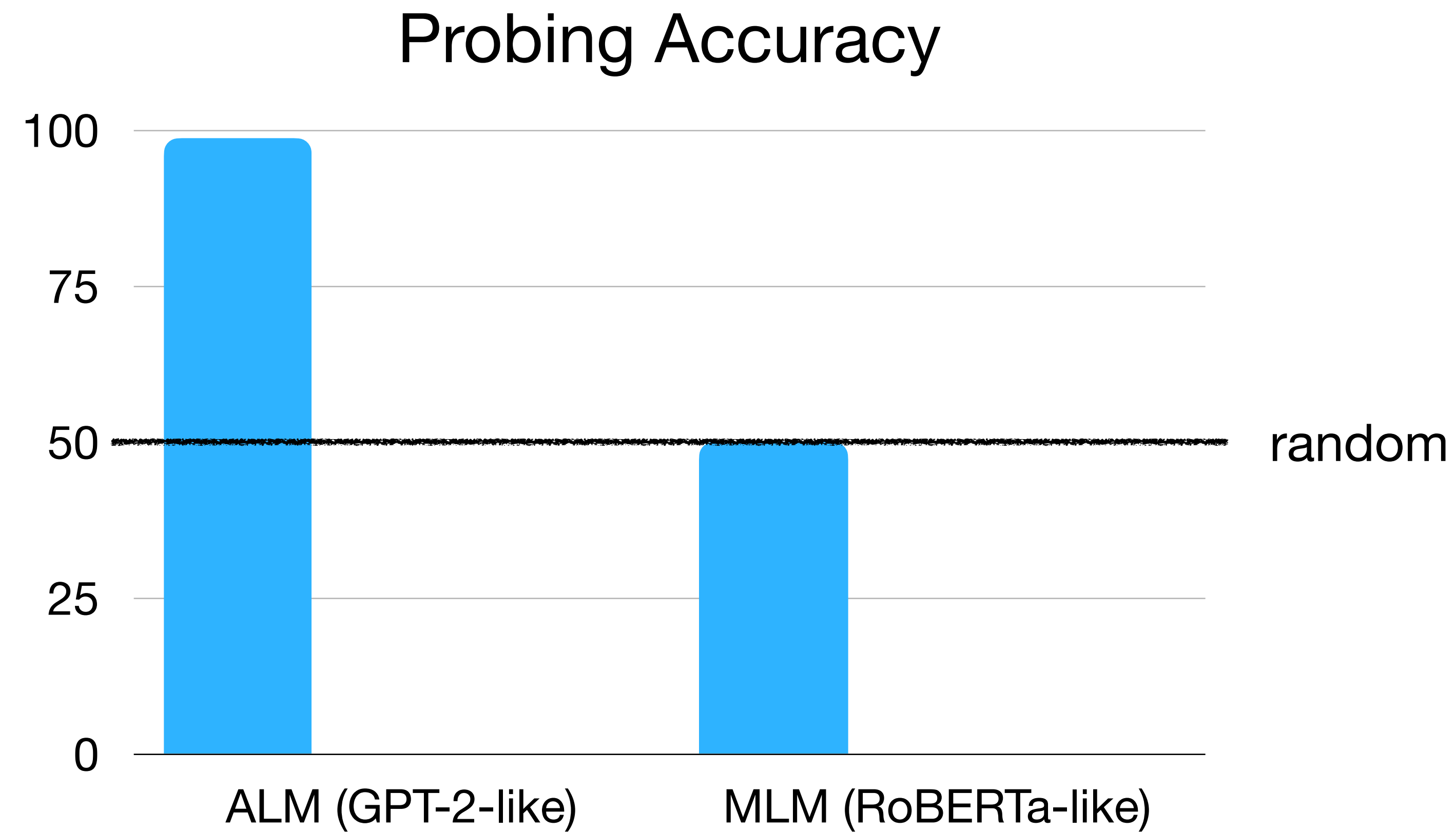
Increasingly more probe parameters

Results



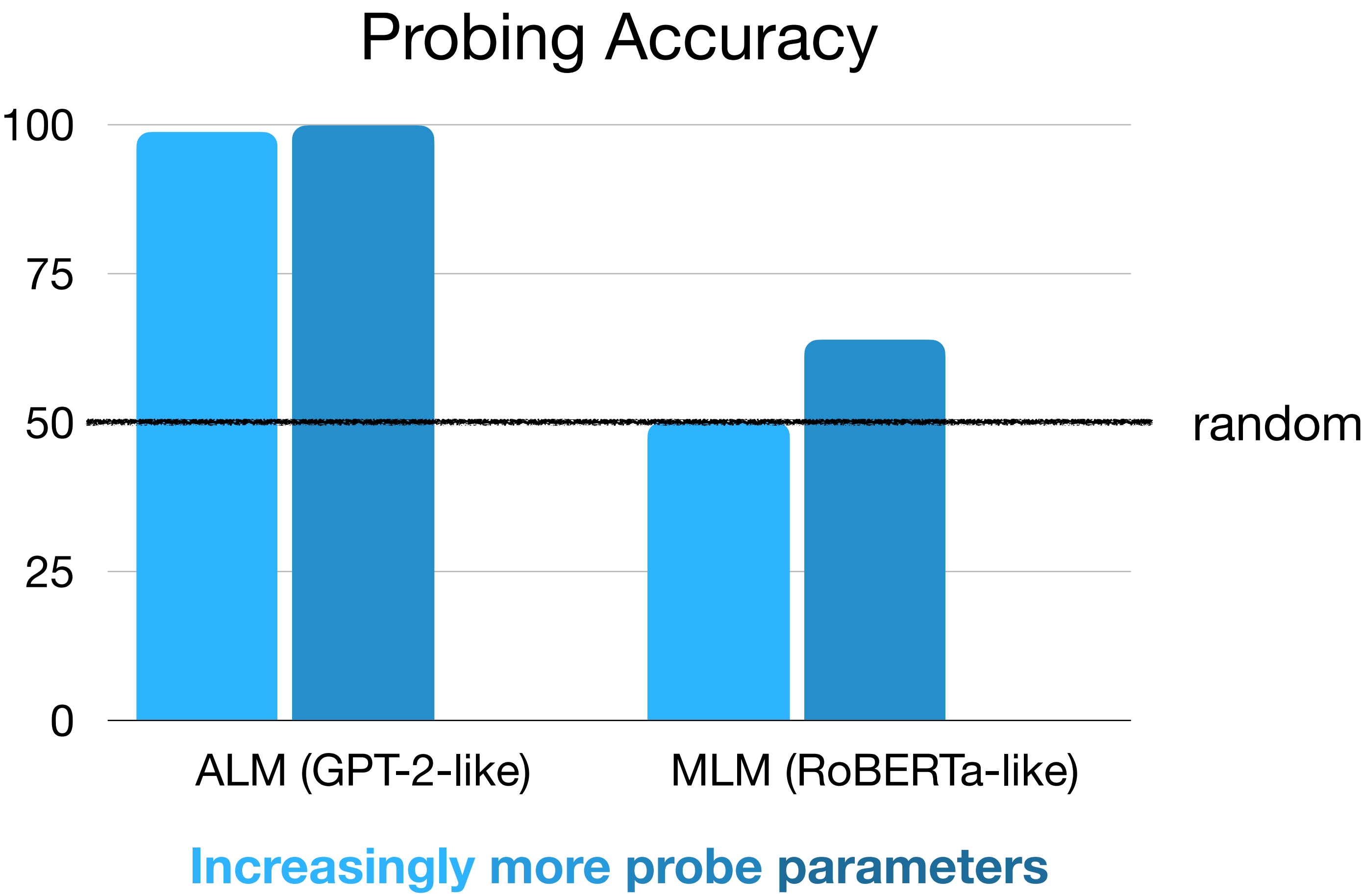
Increasingly more probe parameters

Results

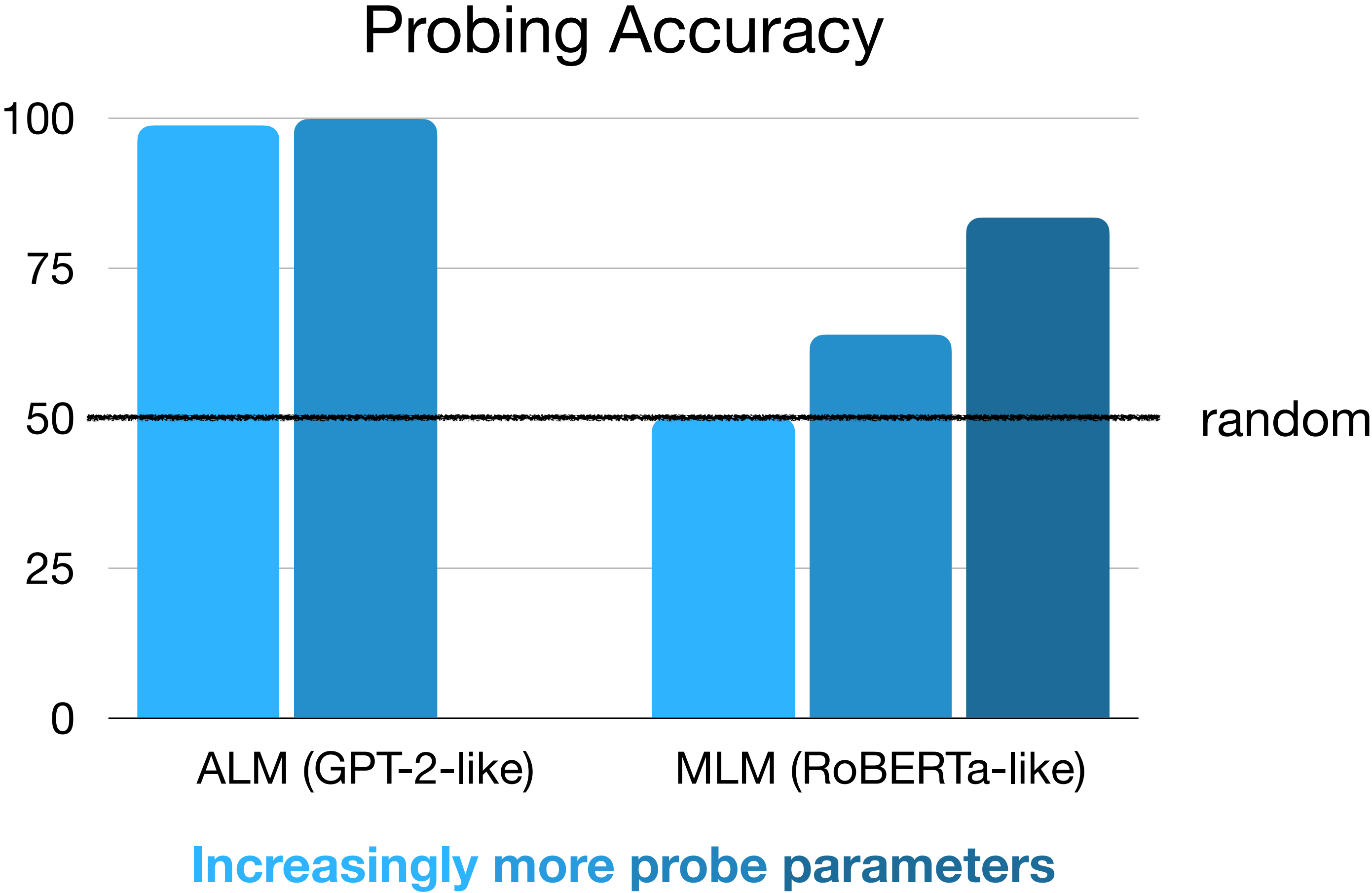


Increasingly more probe parameters

Results



Results



Direct Evaluation

Direct Evaluation

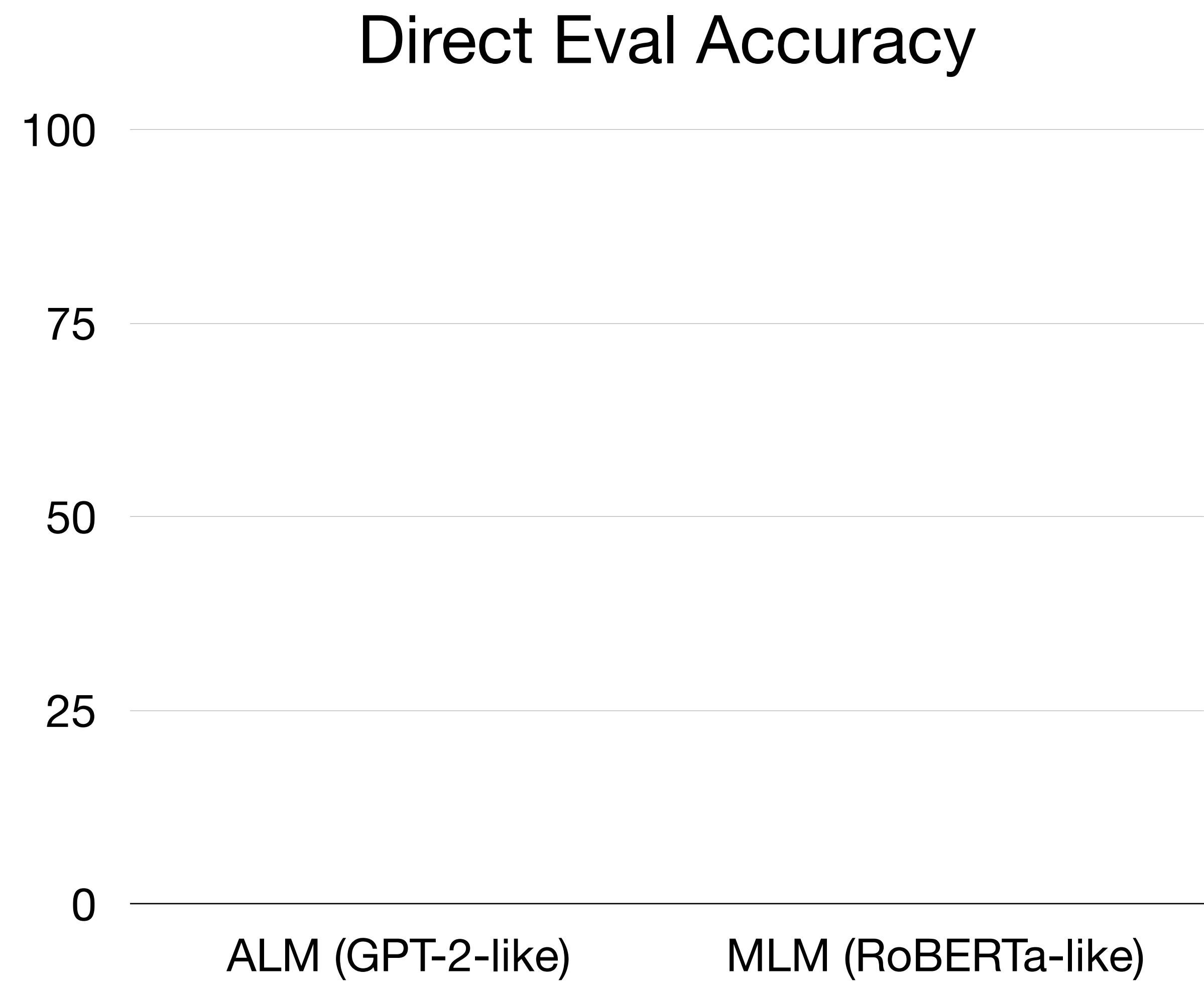
- $((\neg T) \vee F) \vee (\neg T) = \underline{\hspace{1cm}}$

Direct Evaluation

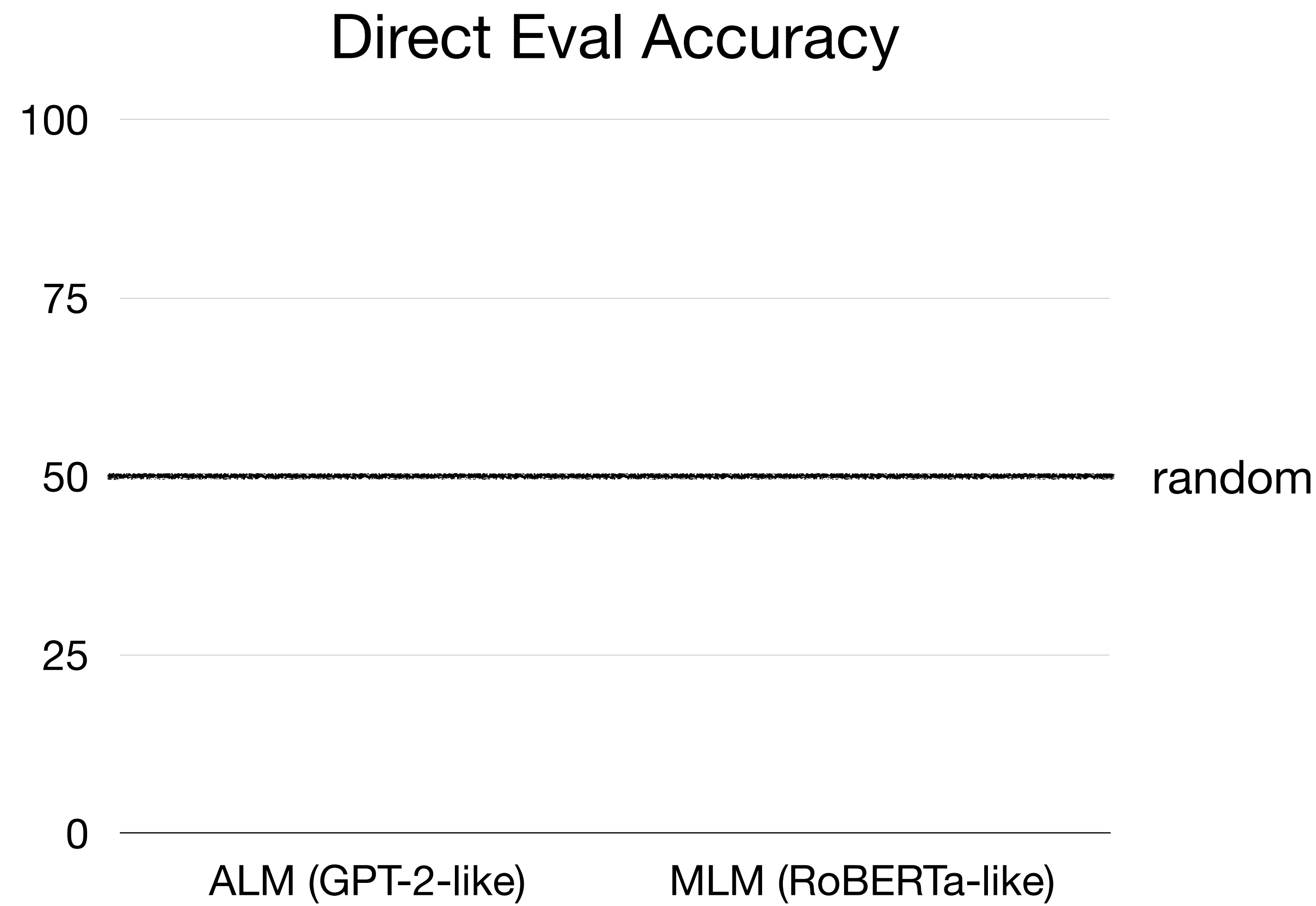
- $((\neg T) \vee F) \vee (\neg T) = \underline{\hspace{1cm}}$
- (small twist, see paper)

Results

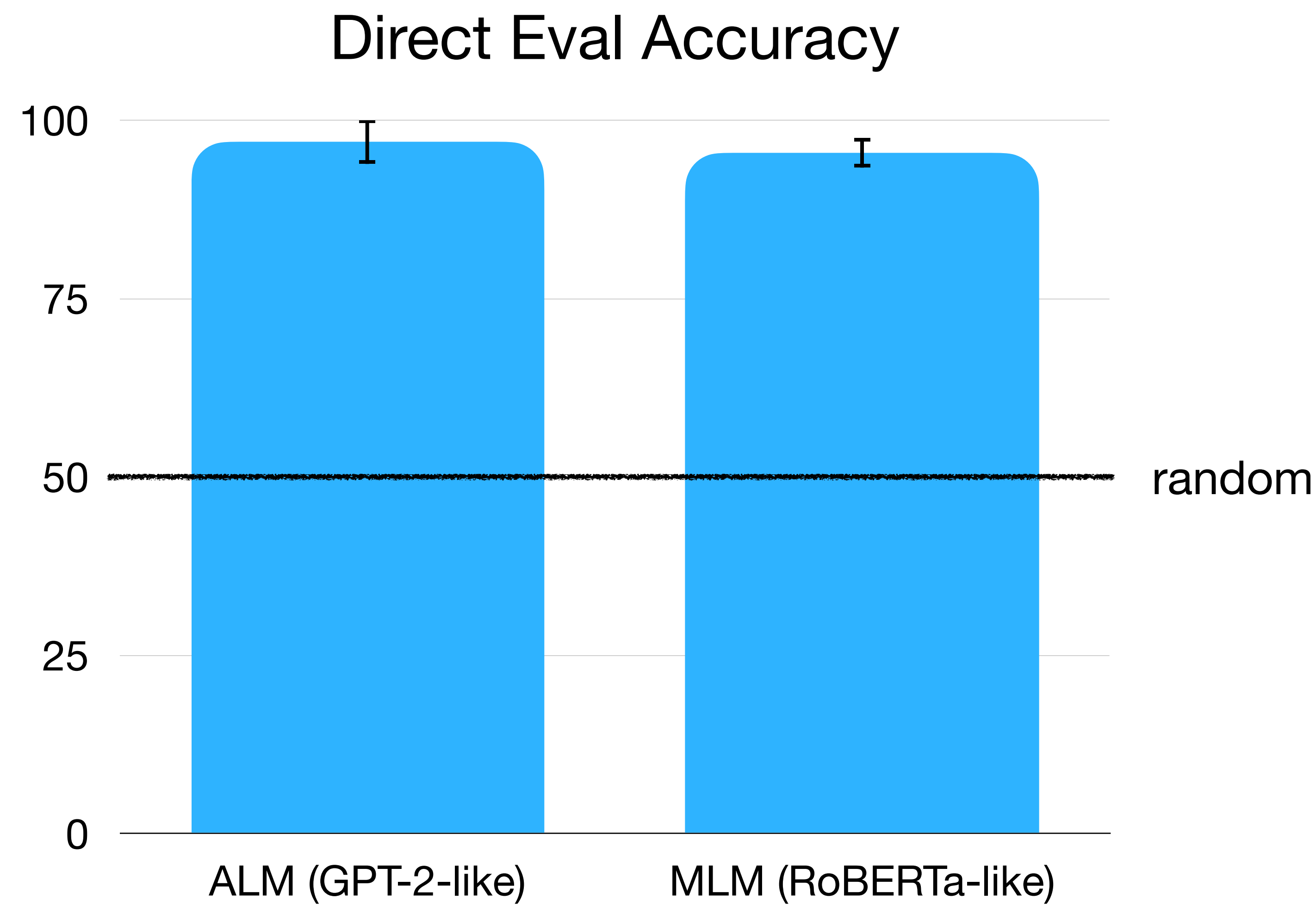
Results



Results



Results



Summary

Summary



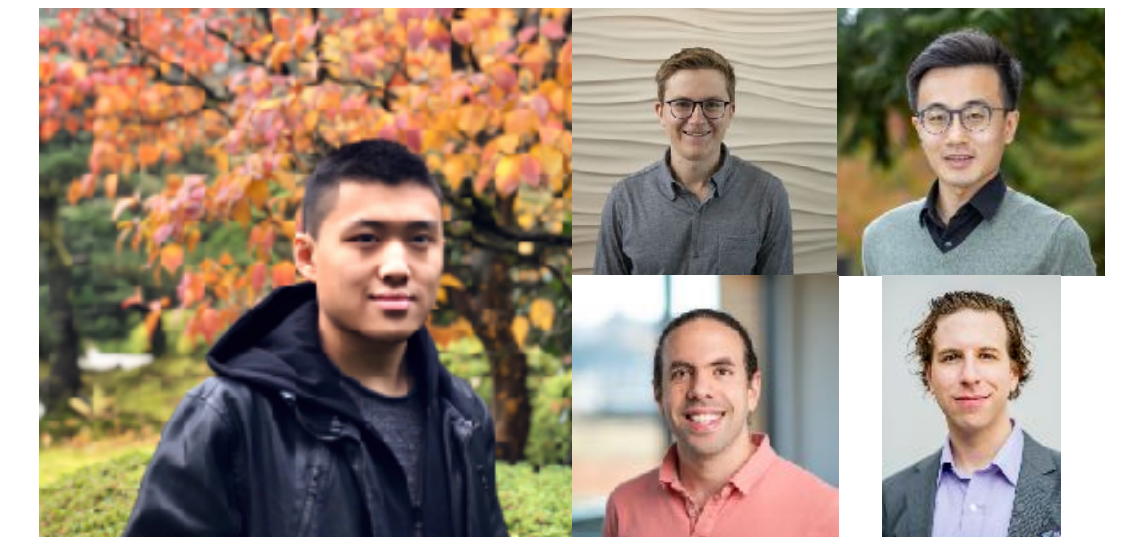
We let GPT-2 complete the simple arithmetic problem *Three plus five equals*. The five responses below [...] show that this problem is beyond the current capability of GPT-2, and, we would argue, any pure LM.

Summary



We let GPT-2 complete the simple arithmetic problem *Three plus five equals*. The five responses below [...] show that this problem is beyond the current capability of GPT-2, and, we would argue, any pure LM.

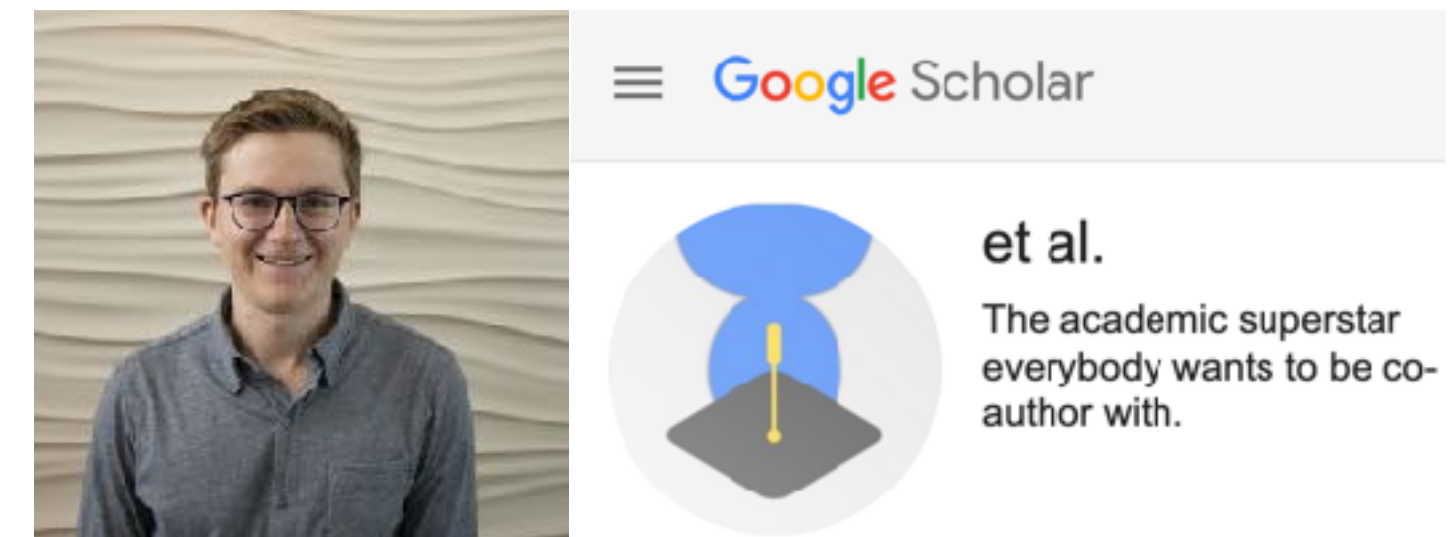
LMs can learn to consistently compare and evaluate the meaning of propositional logic expressions.



What About Other Languages?

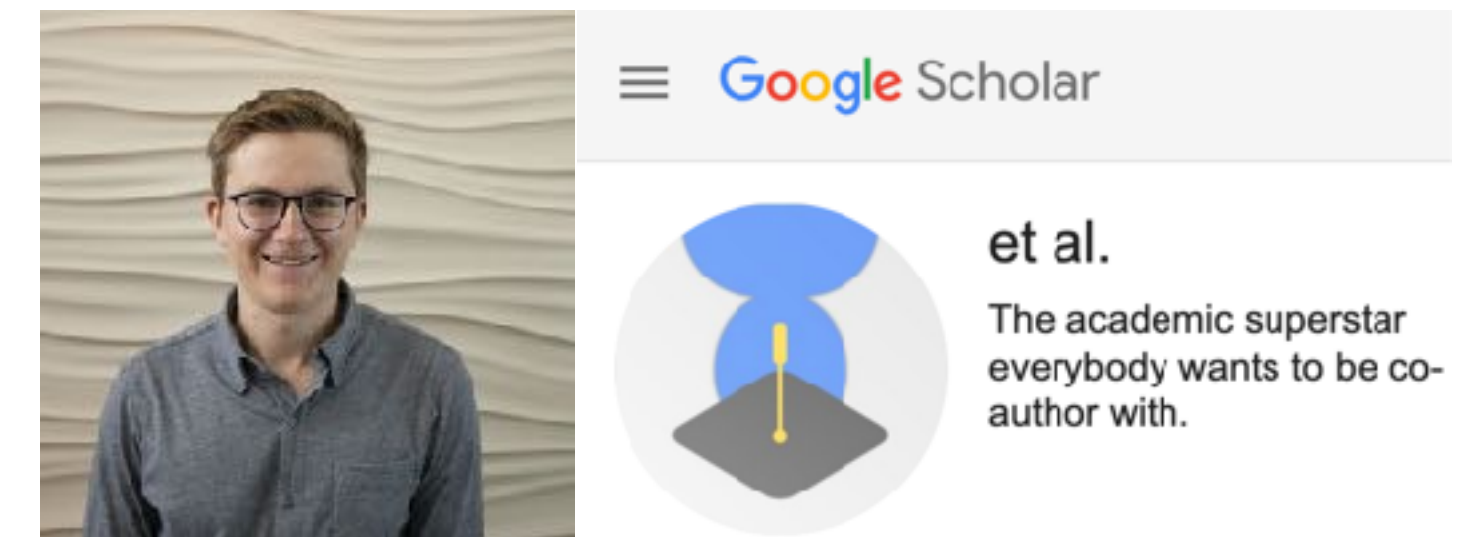
What About Other Languages?

Assertions enable meaning learnability in some languages.



What About Other Languages?

Assertions enable meaning learnability in some languages. ???



Strong Transparency

(i.e., context-independency)

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

$3+5$

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

$3+5$

$x+5$

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

$3+5$

x

$x+5$

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

3+5

x

x+5

date.today()

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

3+5

x

x+5

JUL
17

date.today()

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

$3+5$

Some corgis run.



x

$x+5$

JUL
17

`date.today()`

Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

$3+5$

x

$x+5$

JUL
17

`date.today()`

Some corgis run.

His corgis run.



Strong Transparency

(i.e., context-independency)

- An expression is strongly transparent if its meaning is context-independent
- A language is strongly transparent if all of its expressions are

$((T \wedge (F \vee F)) \vee (T \vee (F \wedge T)))$

$3+5$

x

$x+5$

JUL
17

`date.today()`

Some corgis run.

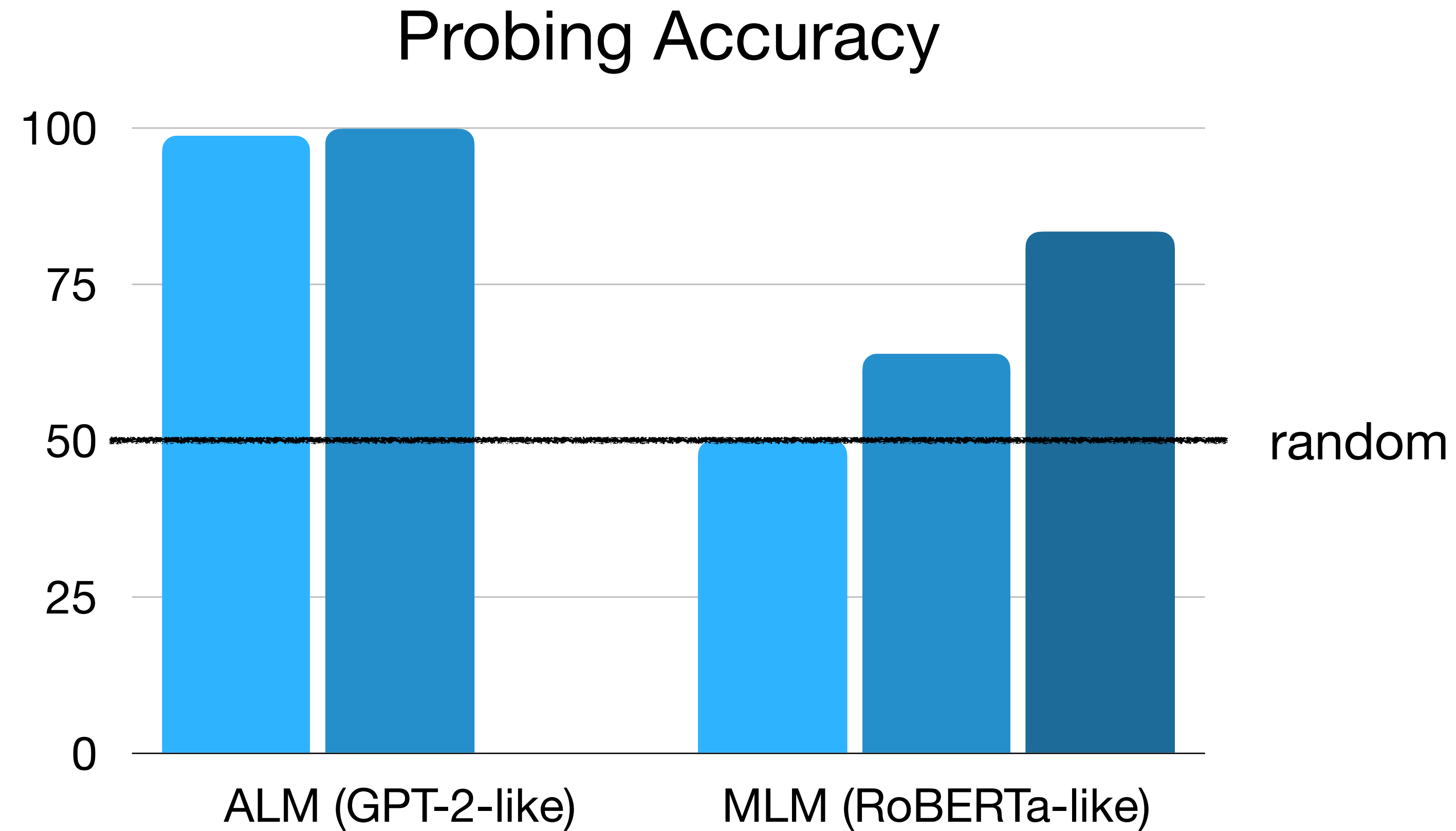
His corgis run.

Today, some corgis ran.



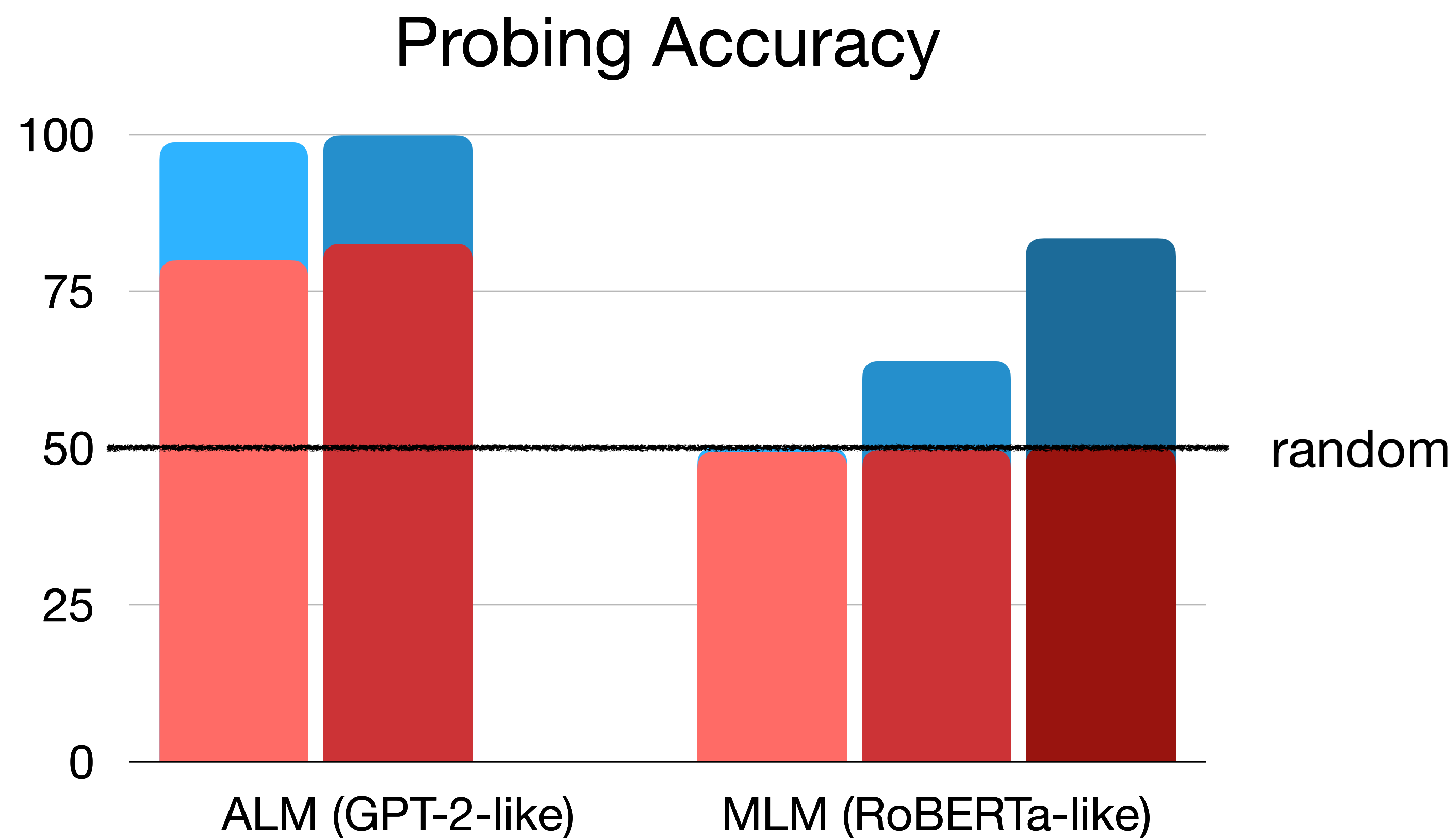
Removing Strong Transparency

Removing Strong Transparency



Increasingly more probe parameters

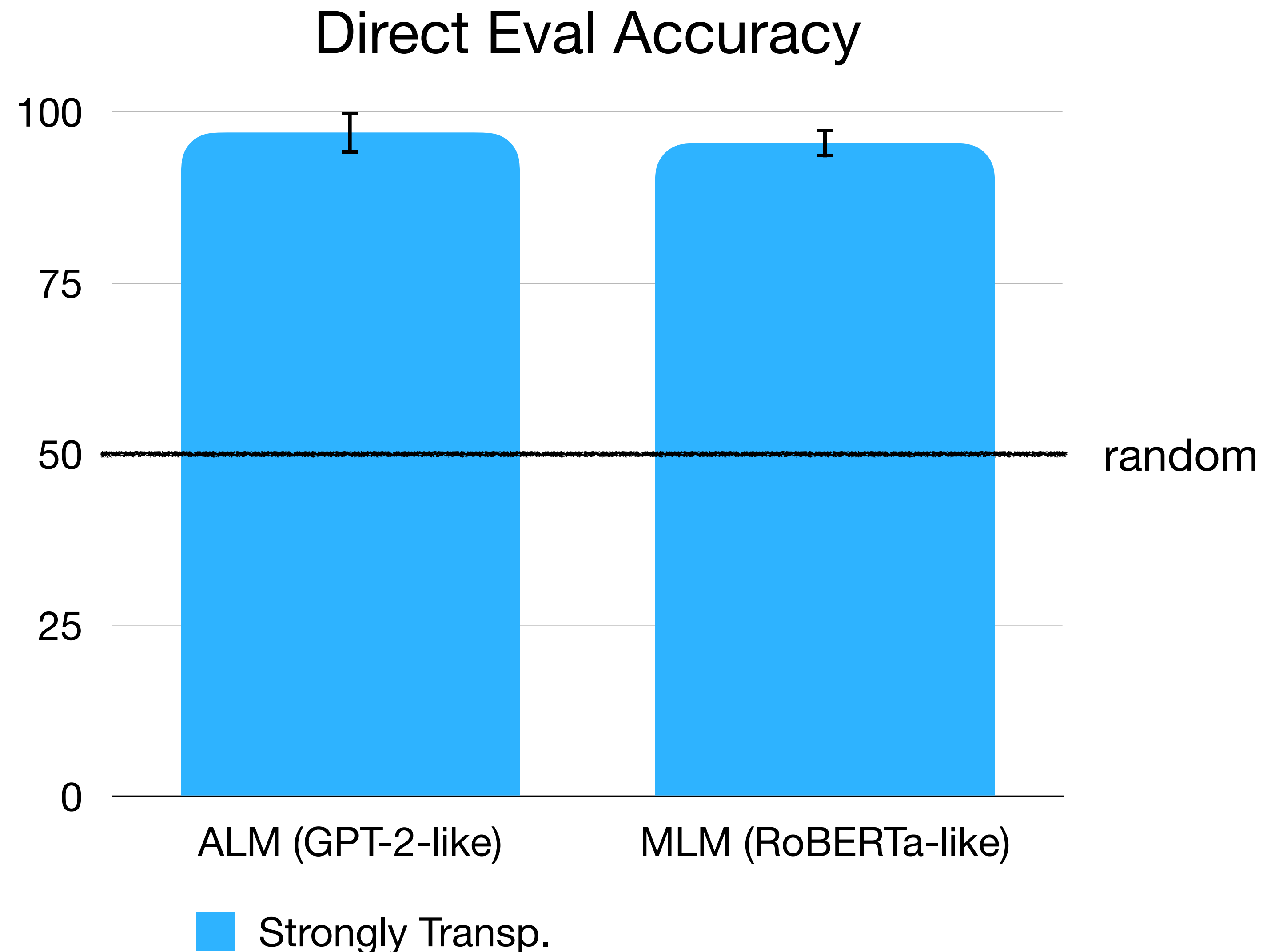
Removing Strong Transparency



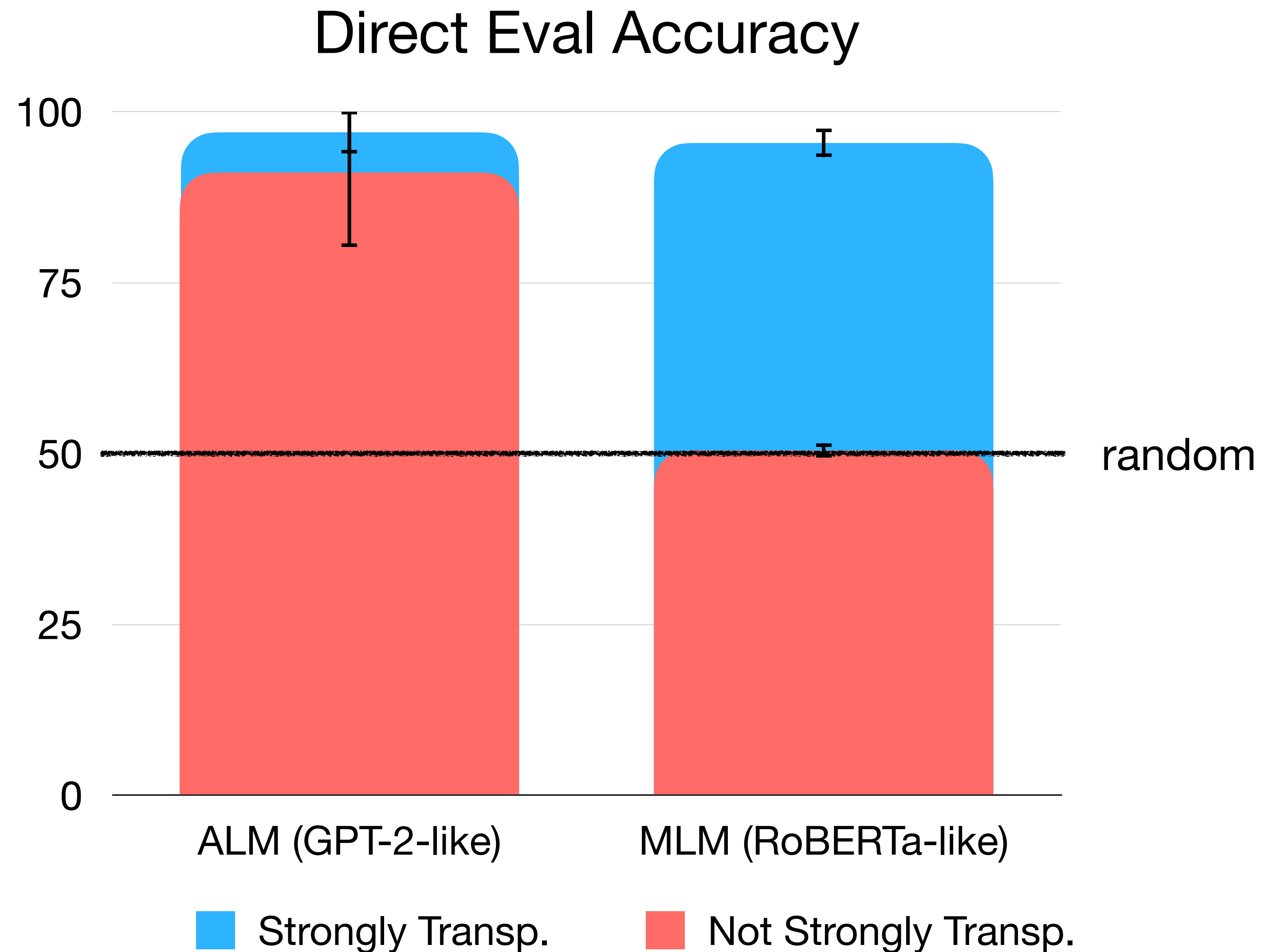
Increasingly more probe parameters

Strongly Transp. Not Strongly Transp.

Removing Strong Transparency

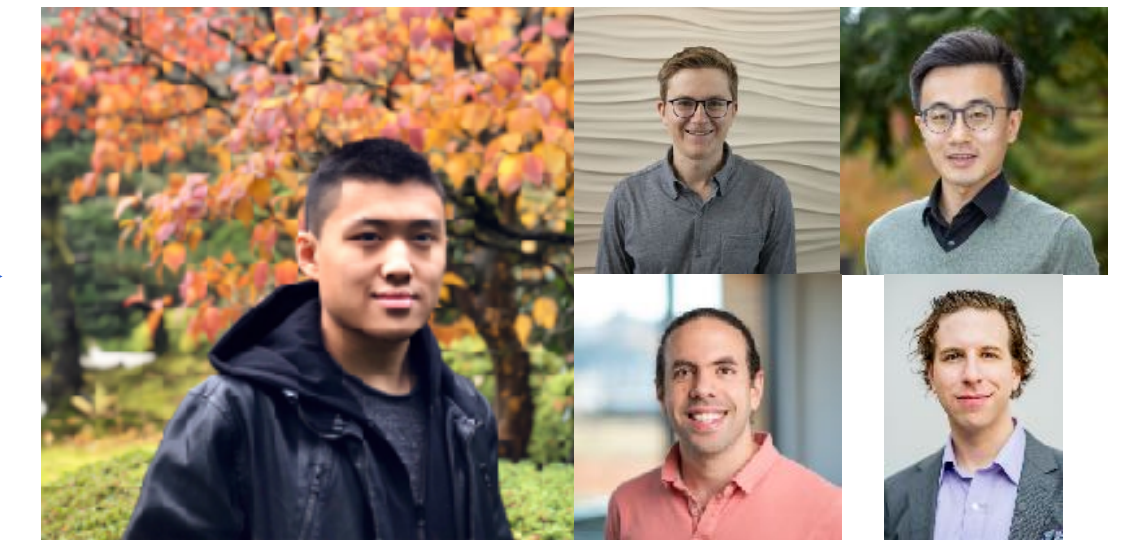


Removing Strong Transparency



Another Summary

LMs can learn the meaning of a strongly transparent language.
And strong transparency is important for this learnability.



But is NL strongly transparent?

Referential Opacity

Foreshadow: it makes NL not strongly transparent

Referential Opacity

Foreshadow: it makes NL not strongly transparent

$[[\text{Superman}]] = [[\text{Clark Kent}]]$

Referential Opacity

Foreshadow: it makes NL not strongly transparent

$[[\text{Superman}]] = [[\text{Clark Kent}]] =$



Referential Opacity

Foreshadow: it makes NL not strongly transparent

[[Superman]] = [[Clark Kent]] =



[[Lois Lane believes Superman is a hero.]]

||

T

Referential Opacity

Foreshadow: it makes NL not strongly transparent

[[Superman]] = [[Clark Kent]] =



[[Lois Lane believes Superman is a hero.]]

||

T

[[Lois Lane believes Clark Kent is a hero.]]

||

F

Referential Opacity

Foreshadow: it makes NL not strongly transparent

$[[\text{Superman}]] = [[\text{Clark Kent}]] =$



$[[\text{Lois Lane believes Superman is a hero.}]] \neq [[\text{Lois Lane believes Clark Kent is a hero.}]]$

\parallel

T

\parallel

F

Referential Opacity

Foreshadow: it makes NL not strongly transparent

$[[\text{Superman}]] = [[\text{Clark Kent}]] =$



propositional attitude verb



$[[\text{Lois Lane believes Superman is a hero.}]] \neq [[\text{Lois Lane believes Clark Kent is a hero.}]]$

\parallel

T

\parallel

F

Formalizing Referential Opacity

Formalizing Referential Opacity

- **Theorem:** A compositional language with referential opacity is not strongly transparent

Formalizing Referential Opacity

- **Theorem:** A compositional language with referential opacity is not strongly transparent
- We know the meaning of strongly transparent languages is learnable

Formalizing Referential Opacity

- **Theorem:** A compositional language with referential opacity is not strongly transparent
- We know the meaning of strongly transparent languages is learnable
- But we saw strong transparency is important for learnability

Formalizing Referential Opacity

- **Theorem:** A compositional language with referential opacity is not strongly transparent
- We know the meaning of strongly transparent languages is learnable
- But we saw strong transparency is important for learnability
- How well do LMs learn this NL phenomenon that is not strongly transparent?

Setup

Setup

- **Data:** $\{(s_1, s_2, y)\}$

Setup

- **Data:** $\{(s_1, s_2, y)\}$
 - She **wants** to meet {Superman/Clark Kent}. $y = \text{Non-equivalent}$

Setup

- **Data:** $\{(s_1, s_2, y)\}$
 - She **wants** to meet {Superman/Clark Kent}. $y = \text{Non-equivalent}$
 - She **managed** to meet {Superman/Clark Kent}. $y = \text{Equivalent}$

Setup

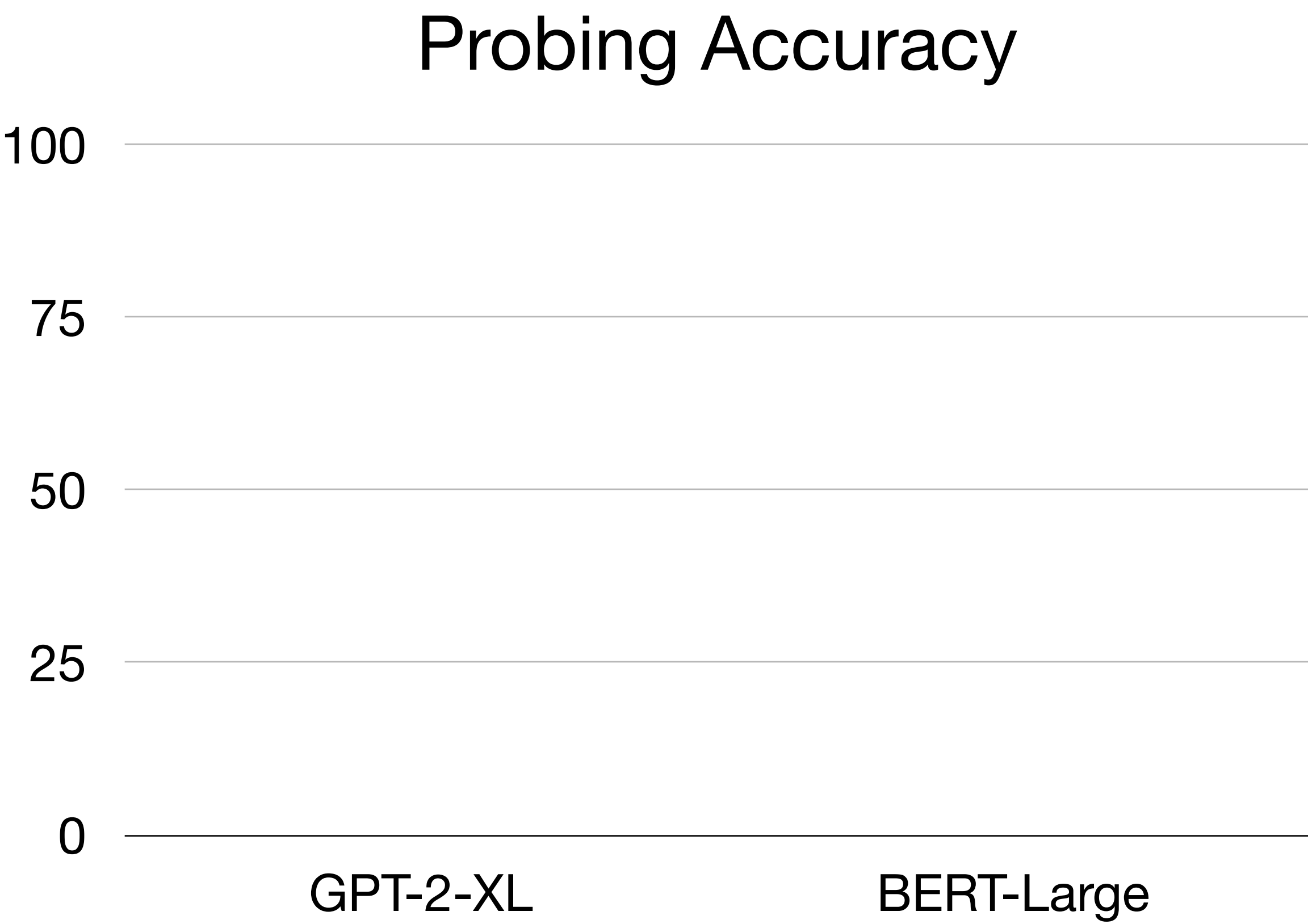
- **Data:** $\{(s_1, s_2, y)\}$
 - She **wants** to meet {Superman/Clark Kent}. $y = \text{Non-equivalent}$
 - She **managed** to meet {Superman/Clark Kent}. $y = \text{Equivalent}$
- **Models:** pretrained GPT-2-XL, BERT-large

Setup

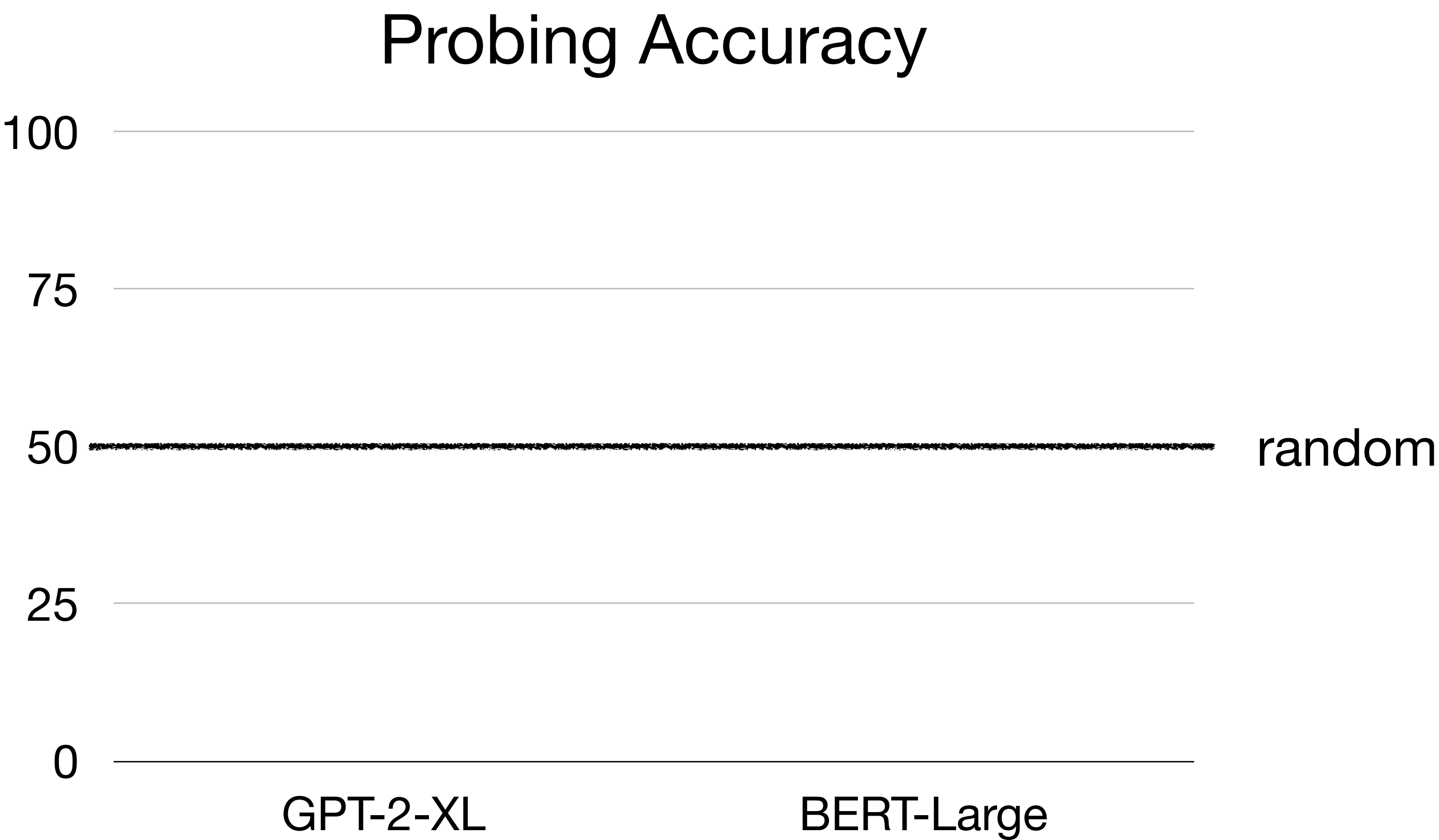
- **Data:** $\{(s_1, s_2, y)\}$
 - She **wants** to meet {Superman/Clark Kent}. $y = \text{Non-equivalent}$
 - She **managed** to meet {Superman/Clark Kent}. $y = \text{Equivalent}$
- **Models:** pretrained GPT-2-XL, BERT-large
- **Methods:** probing and similarity-based analysis

Probing Results

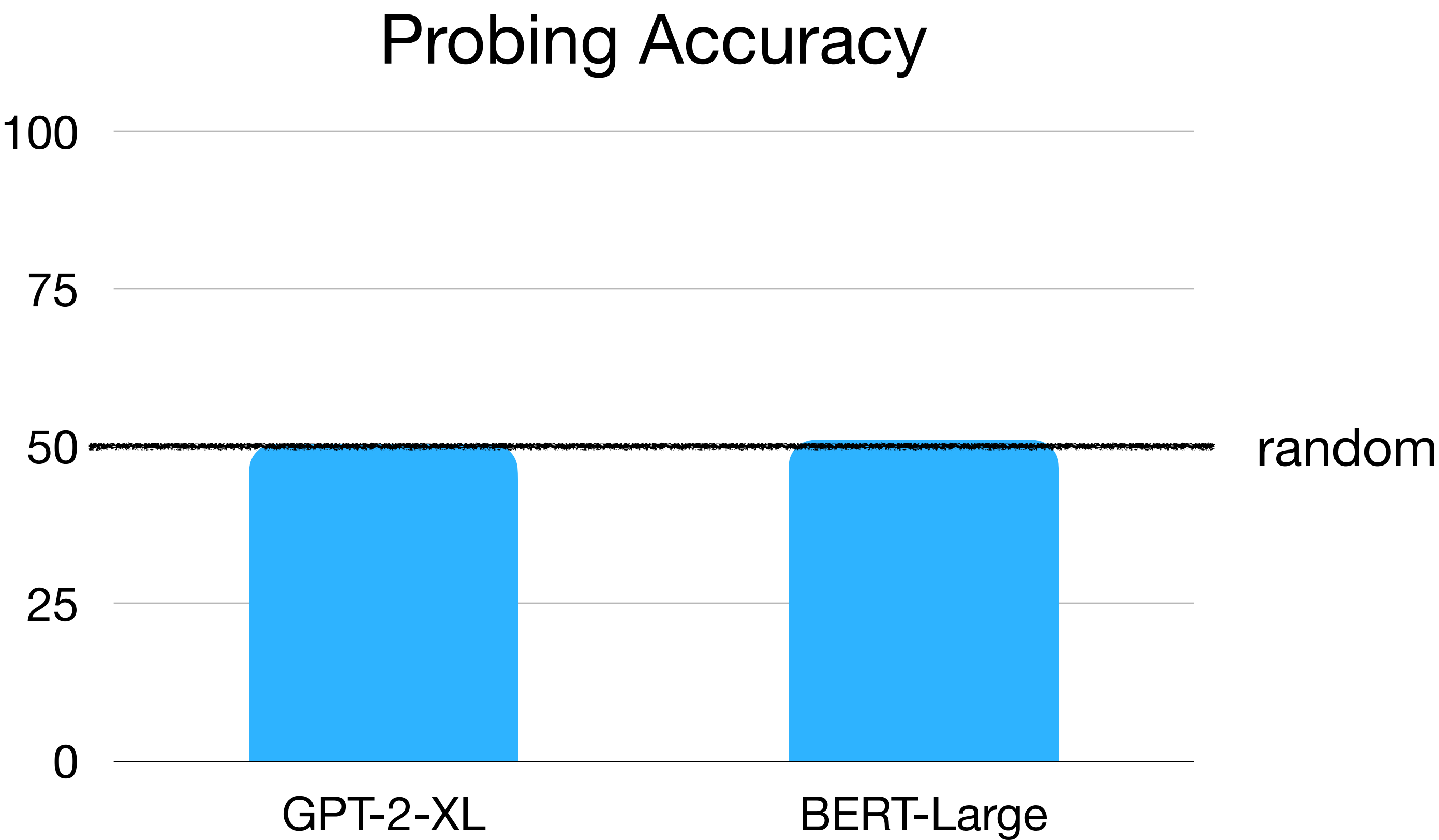
Probing Results



Probing Results

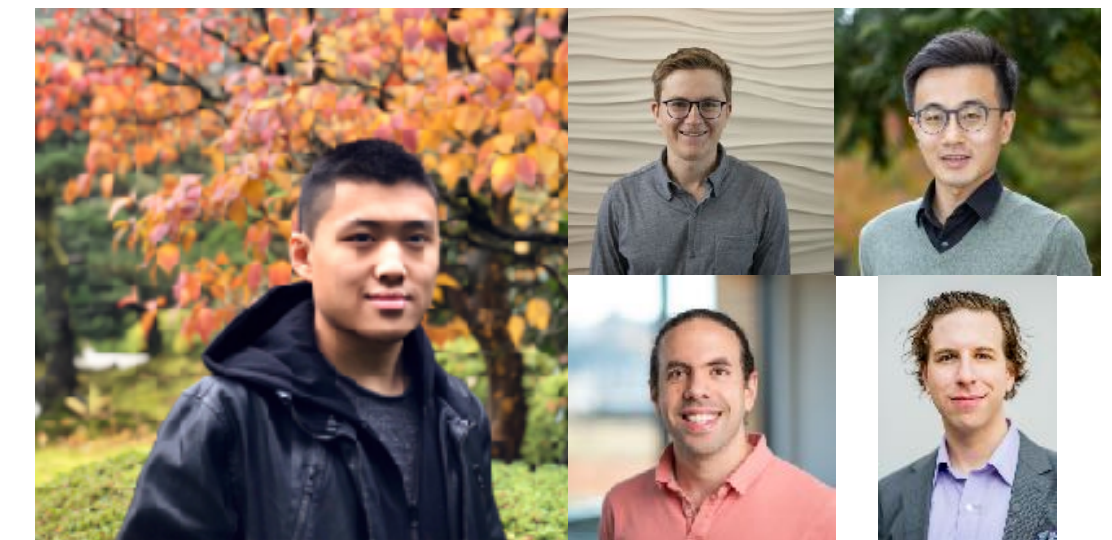


Probing Results



Yet Another Summary

Although LMs could learn the meaning of a strongly transparent language, they don't well-represent referential opacity and hence the meaning of the entirety of NL.



Conclusions

Conclusions

- Aligning with the theoretical guarantee, current LM architectures & objectives can learn the meaning of a strongly transparent language

Conclusions

- Aligning with the theoretical guarantee, current LM architectures & objectives can learn the meaning of a strongly transparent language
- Strong transparency plays a big part in this learnability

Conclusions

- Aligning with the theoretical guarantee, current LM architectures & objectives can learn the meaning of a strongly transparent language
- Strong transparency plays a big part in this learnability
 - Though learnability is not completely destroyed w/o strong transparency

Conclusions

- Aligning with the theoretical guarantee, current LM architectures & objectives can learn the meaning of a strongly transparent language
- Strong transparency plays a big part in this learnability
 - Though learnability is not completely destroyed w/o strong transparency
- On NL, there is no evidence at all of LMs representing referential opacity, a phenomenon that is not strongly transparent

Propositional Logic vs. NL

Propositional Logic vs. NL

- Why did we see >random probing/eval accuracy on the perturbed propositional logic, but not referential opacity?

Propositional Logic vs. NL

- Why did we see >random probing/eval accuracy on the perturbed propositional logic, but not referential opacity?
 - Maybe referential opacity is just harder

Propositional Logic vs. NL

- Why did we see >random probing/eval accuracy on the perturbed propositional logic, but not referential opacity?
 - Maybe referential opacity is just harder
 - Maybe it's because of the large variation in NL, with sentences that are untruthful, subjective, etc.

Propositional Logic vs. NL

- Why did we see >random probing/eval accuracy on the perturbed propositional logic, but not referential opacity?
 - Maybe referential opacity is just harder
 - Maybe it's because of the large variation in NL, with sentences that are untruthful, subjective, etc.
 - Or maybe...

Encore: Grounding =

Encore: Grounding =

$((\neg T) \wedge (\neg(T \vee (\neg F)))) = (T \vee (\neg(\neg((\neg T) \vee (\neg(\neg F)))))$
 $(\neg(\neg(\neg((F \wedge ((F \wedge F) \wedge F)) \wedge F) \wedge (\neg T)))) = ((T \wedge T) \wedge ((\neg F) \vee (\neg F)))$
 $((\neg(\neg(\neg(\neg(\neg T)))) \vee T) \vee T) \wedge (\neg(\neg T)) = ((\neg F) \vee (\neg(T \wedge (T \vee T))))$
 $((T \wedge (F \vee F)) \vee (T \vee (F \wedge T))) = (\neg((\neg T) \wedge (\neg(\neg(\neg(\neg F)) \vee F)) \vee (T \wedge T)))$
 $((\neg(\neg F)) \wedge (\neg F)) \wedge ((\neg F) \vee F) \wedge F) = (F \wedge (\neg(\neg(F \vee (\neg(F \vee (\neg T))) \wedge T))))$
 $((T \vee (\neg(T \wedge (T \vee (\neg(F \vee (\neg F))))) \vee T) = (\neg(\neg(T \vee (\neg(\neg(\neg(T \wedge F))))) \wedge F))$
 $(\neg(\neg(\neg F)) \vee ((\neg F) \vee (T \vee (\neg(T \vee T))))) = ((F \wedge (\neg T)) \wedge ((\neg F) \wedge F)) \wedge F$
 $(F \wedge (F \wedge (\neg(\neg(T \vee T)) \wedge (\neg T)))) = (\neg(\neg((\neg T) \vee F) \vee (F \vee (\neg(\neg F)))))$
 $(F \wedge (F \wedge (\neg((F \vee F) \vee (\neg(\neg T))))) = (\neg(((\neg(T \wedge T)) \vee (\neg F)) \vee (\neg T)) \wedge (\neg F))$

Encore: Grounding =

$((\neg T) \wedge (\neg (Tv(\neg F)))) = (Tv(\neg(\neg((\neg T) \vee (\neg(\neg F)))))$
 $(\neg(\neg(\neg((F \wedge ((F \wedge F) \wedge F)) \wedge F) \wedge (\neg T)))) = ((T \wedge T) \wedge ((\neg F) \vee (\neg F)))$
 $((\neg(\neg(\neg(\neg(\neg T)))) \vee T) \vee T) \wedge (\neg(\neg T)) = ((\neg F) \vee (\neg(T \wedge (TvT))))$
 $((T \wedge (F \vee F)) \vee (Tv(F \wedge T))) = (\neg((\neg T) \wedge (\neg(\neg(\neg(\neg F)) \vee F)) \vee (T \wedge T)))$
 $((\neg(\neg F)) \wedge (\neg F)) \wedge ((\neg F) \vee F) \wedge F = (F \wedge (\neg(\neg(F \vee (\neg(F \vee (\neg T))) \wedge T))))$
 $((Tv(\neg(T \wedge (Tv(\neg(F \vee (\neg F))))) \vee T) = (\neg((\neg(Tv(\neg(\neg(\neg(T \wedge F))))) \wedge F))$
 $(\neg((\neg(\neg F)) \vee ((\neg F) \vee (Tv(\neg(TvT))))) = ((F \wedge (\neg T)) \wedge ((\neg F) \wedge F)) \wedge F$
 $(F \wedge (F \wedge (\neg((\neg(TvT)) \wedge (\neg T)))) = (\neg(\neg((\neg T) \vee F) \vee (F \vee (\neg(\neg F)))))$
 $(F \wedge (F \wedge (\neg((F \vee F) \vee (\neg(\neg T))))) = (\neg(((\neg(T \wedge T)) \vee (\neg F)) \vee (\neg T)) \wedge (\neg F))$

Probing Accuracy

a=b
50.5

Encore: Grounding =

$((\neg T) \wedge (\neg(T \vee (\neg F)))) = (T \vee (\neg(\neg((\neg T) \vee (\neg(\neg F)))))$
 $(\neg(\neg(\neg((F \wedge ((F \wedge F) \wedge F)) \wedge F) \wedge (\neg T)))) = ((T \wedge T) \wedge ((\neg F) \vee (\neg F)))$
 $((\neg(\neg(\neg(\neg(\neg T)))) \vee T) \vee T) \wedge (\neg(\neg T)) = ((\neg F) \vee (\neg(T \wedge (T \vee T))))$
 $((T \wedge (F \vee F)) \vee (T \vee (F \wedge T))) = (\neg((\neg T) \wedge (\neg(\neg(\neg(\neg F)) \vee F)) \vee (T \wedge T)))$
 $((\neg(\neg F)) \wedge (\neg F)) \wedge ((\neg F) \vee F) \wedge F = (F \wedge (\neg(\neg(F \vee (\neg(F \vee (\neg T))) \wedge T))))$
 $((T \vee (\neg(T \wedge (T \vee (\neg(F \vee (\neg F))))) \vee T) = (\neg((\neg(T \vee (\neg(\neg(\neg(T \wedge F))))) \wedge F))$
 $(\neg((\neg(\neg F)) \vee ((\neg F) \vee (T \vee (\neg(T \vee T))))) = ((F \wedge (\neg T)) \wedge ((\neg F) \wedge F) \wedge F)$
 $(F \wedge (F \wedge (\neg((\neg(T \vee T)) \wedge (\neg T)))) = (\neg(\neg((\neg T) \vee F) \vee (F \vee (\neg(\neg F)))))$
 $(F \wedge (F \wedge (\neg((F \vee F) \vee (\neg(\neg T))))) = (\neg(((\neg(T \wedge T)) \vee (\neg F)) \vee (\neg T)) \wedge (\neg F))$

Probing Accuracy

-Reflexivity	+Reflexivity
a=b 50.5	a=b, a=a, b=b 92.7

Encore: Grounding =

$((\neg T) \wedge (\neg(T \vee (\neg F)))) = (T \vee (\neg(\neg((\neg T) \vee (\neg(\neg F)))))$
 $(\neg(\neg(\neg((F \wedge ((F \wedge F) \wedge F)) \wedge F) \wedge (\neg T)))) = ((T \wedge T) \wedge ((\neg F) \vee (\neg F)))$
 $((\neg(\neg(\neg(\neg(\neg T)))) \vee T) \vee T) \wedge (\neg(\neg T)) = ((\neg F) \vee (\neg(T \wedge (T \vee T))))$
 $((T \wedge (F \vee F)) \vee (T \vee (F \wedge T))) = (\neg((\neg T) \wedge (\neg(\neg(\neg(\neg F)) \vee F)) \vee (T \wedge T)))$
 $((\neg(\neg F)) \wedge (\neg F)) \wedge ((\neg F) \vee F) \wedge F = (F \wedge (\neg(\neg(F \vee (\neg(F \vee (\neg T)) \wedge T))))$
 $((T \vee (\neg(T \wedge (T \vee (\neg(F \vee (\neg F))))) \vee T) = (\neg(\neg(T \vee (\neg(\neg(\neg(T \wedge F))))) \wedge F))$
 $(\neg(\neg(\neg F)) \vee ((\neg F) \vee (T \vee (\neg(T \vee T)))) = ((F \wedge (\neg T)) \wedge ((\neg F) \wedge F)) \wedge F$
 $(F \wedge (F \wedge (\neg(\neg(T \vee T)) \wedge (\neg T)))) = (\neg(\neg((\neg T) \vee F) \vee (F \vee (\neg(\neg F)))))$
 $(F \wedge (F \wedge (\neg((F \vee F) \vee (\neg(\neg T))))) = (\neg(((\neg(T \wedge T)) \vee (\neg F)) \vee (\neg T)) \wedge (\neg F))$

Probing Accuracy

	-Reflexivity	+Reflexivity
-Symmetry	a=b 50.5	a=b, a=a, b=b 92.7
+Symmetry	a=b, b=a 50.3	a=b, b=a, a=a, b=b 98.8

Propositional Logic vs. NL

- Why did we see >random probing accuracy on the perturbed propositional logic, but not referential opacity?
 - Maybe referential opacity is just harder
 - Maybe it's because of the variation in NL, with sentences that are untruthful, subjective, etc.
 - Or maybe...

Propositional Logic vs. NL

- Why did we see >random probing accuracy on the perturbed propositional logic, but not referential opacity?
 - Maybe referential opacity is just harder
 - Maybe it's because of the variation in NL, with sentences that are untruthful, subjective, etc.
 - Or maybe...
 - We don't have such an explicit representation of equivalence in NL pretraining

- Aligning with the theoretical guarantee, current LM architectures & objectives can learn the meaning of a strongly transparent language
- Strong transparency plays a big part in this learnability
 - Though learnability is not completely destroyed w/o strong transparency
- On NL, there is no evidence at all of LMs representing referential opacity, a phenomenon that is not strongly transparent
- Careful design of the pretraining data/setup is crucial