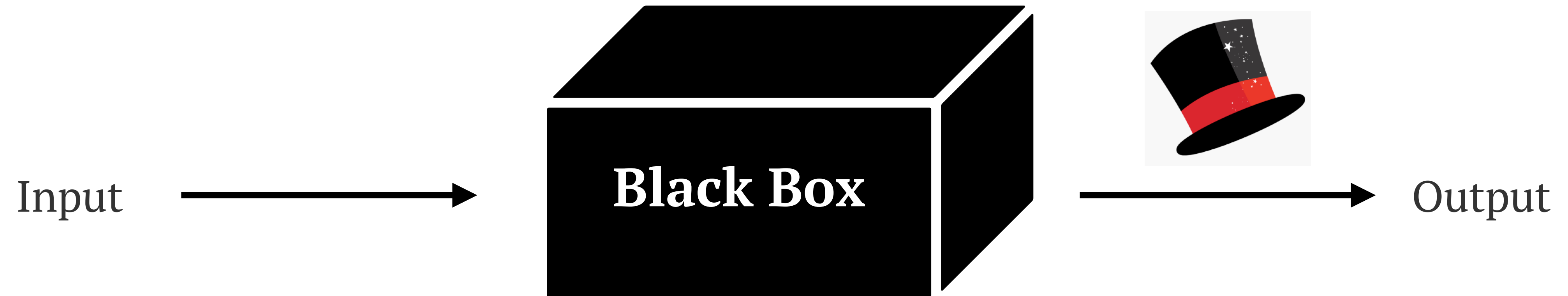# Model Interpretability

## COMP3361 — Week 12

Lingpeng Kong

Department of Computer Science, The University of Hong Kong

Many materials from Zhiyong Wu with Special Thanks!

# Black Box Models

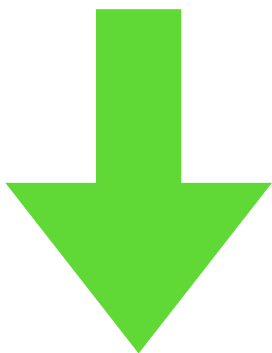Input $\rightarrow$ **Black Box** $\rightarrow$ Output

End-to-End Models
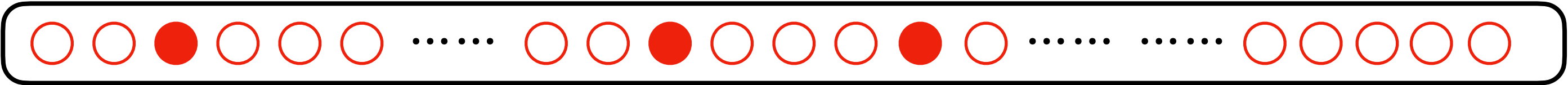
# Linear Models in NLP

A Great movie ! I really like the cast .

Words:

featurized

A
great
movie
!
I
really
like
its
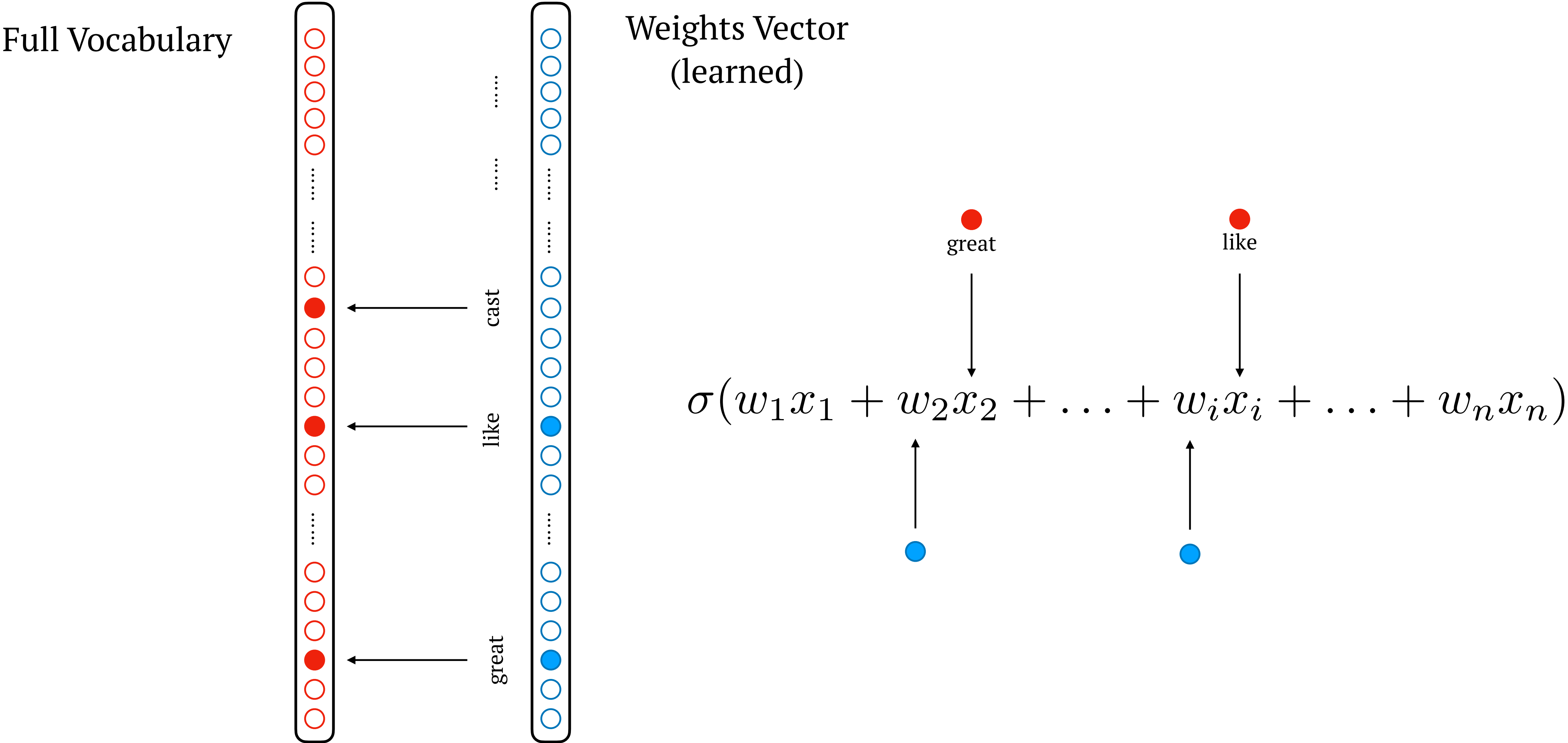cast
.

Full Vocabulary

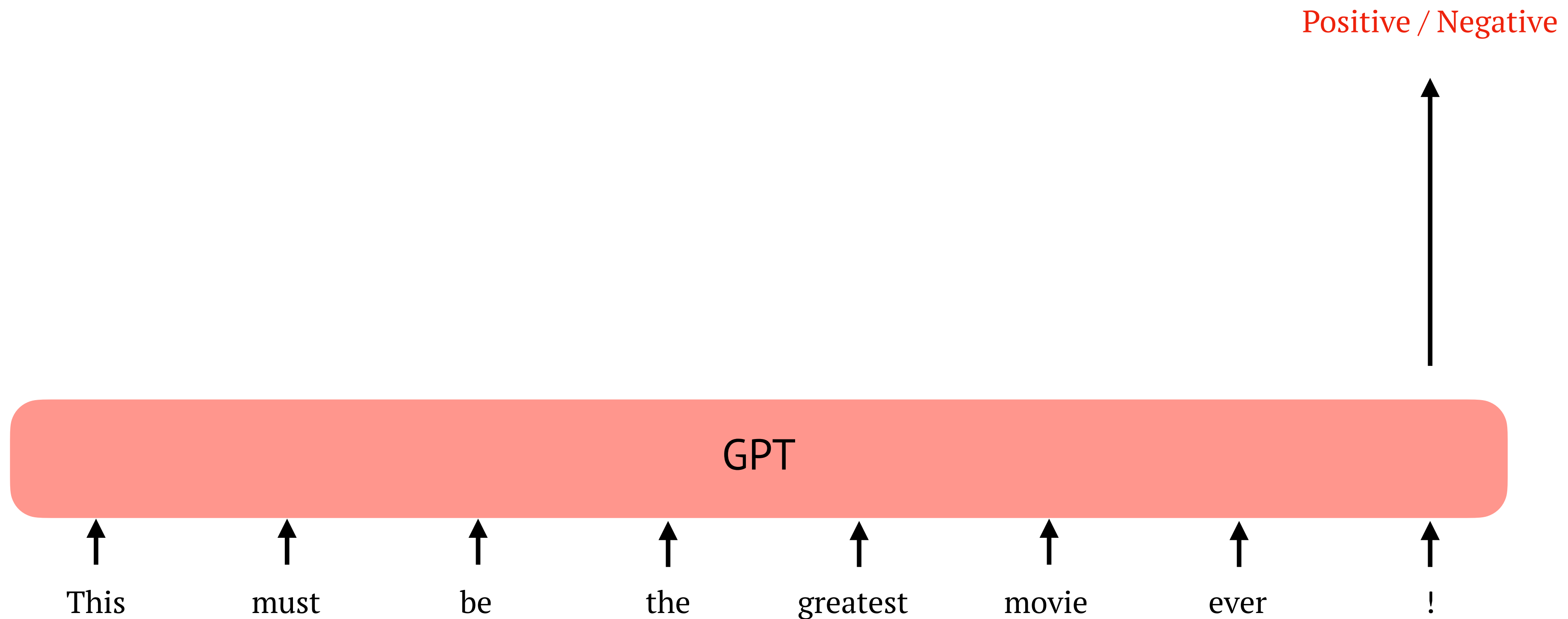great       like    cast     ......    ......

Weights Vector
(learned)

# Linear Models in NLP

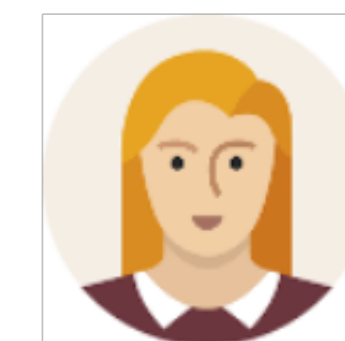Full Vocabulary

Weights Vector
(learned)

$$\sigma(w_1 x_1 + w_2 x_2 + \ldots + w_i x_i + \ldots + w_n x_n)$$

great

like

cast

like

great

# Neural Models in NLP

Positive / Negative

GPT

This    must    be    the    greatest    movie    ever    !

# Neural Models in NLP

# The Great Sucess of Large-scale Language Models

# A Hypothesis



| Linguistic Knowledge | Factual Knowledge |
|---|---|

Internalize → Google BERT ← Internalize

# Probing

# Probing

$$\mathrm{arc\_score}\ (\ \bigcirc,\bigcirc\ )\ =\ \boxed{\text{PROB}}\ \bigcirc + \boxed{\text{PROB}}\ \bigcirc + \boxed{\text{PROB}}$$



BERT

**Freeze**

# Model Probing



Shall we give credit to the **representation**? and/or the **probe**?

# Parameter-free Probing

**Perturbed Masking**



(Wu et al., 2020)

# Parameter-free Probing

Example: Calculate impact *sit* has on *Cats*:  $f(Cats|sit) = d(e, e')$

$e = E(Cats|S\backslash\{Cats\})$                    $e' = E(Cats|S\backslash\{Cats, sit\})$



(Wu et al., 2020)

# Parameter-free Probing

| | Cats | sit | on | the | mat |
|---|---|---|---|---|---|
| **Cats** | - | f(Cat,sit) | f(Cat,on) | f(Cat,the) | f(Cat,mat) |
| **sit** | f(sit, Cats) | - | ... | ... | ... |
| **on** | f(on, Cats) | ... | - | ... | ... |
| **the** | f(the, Cats) | ... | ... | - | ... |
| **mat** | f(mat, Cats) | ... | ... | ... | - |

Supervised Probe: learning to map representations to task
Here: Impact Matrix + task specific algorithm => task

(Wu et al., 2020)

# Dependency Knowledge



Quality of the extracted tree structure on WSJ10 (UAS score)

Would BERT learn **better** dependency structures?

(Wu et al., 2020)

# Empirical Evaluation of the BERT Syntax

Aspect-level sentiment classification:
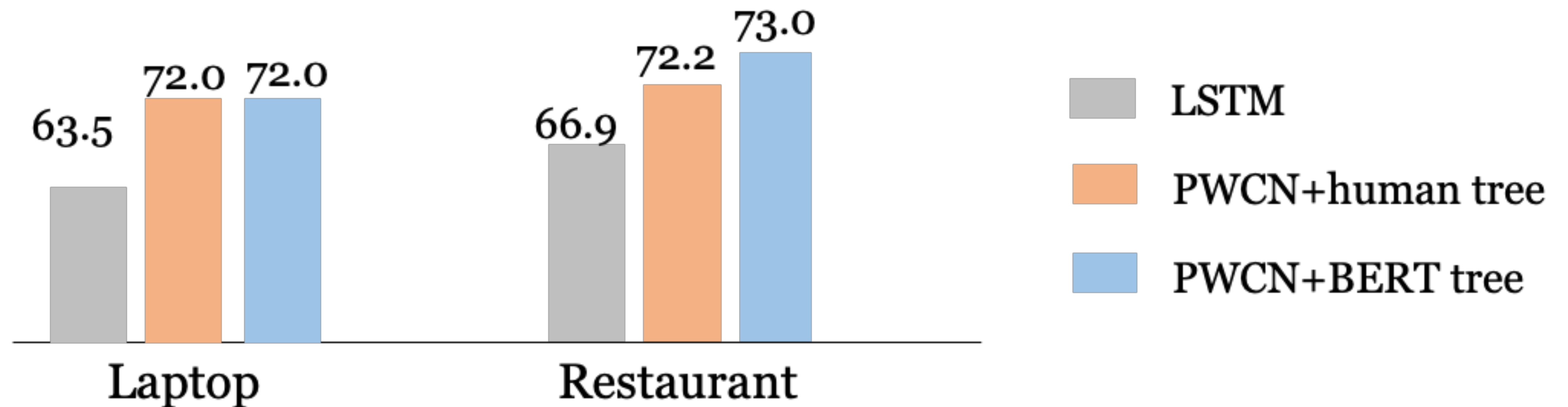    Their food, in my opinion, is ok, but the service is terrible.
Input sentence: **s**
Parser generated dep tree for s: **human tree (linguists-defined syntax)**
BERT generated dep tree for s: **BERT tree (BERT syntax)**



(Wu et al., 2020)

# Empirical Evaluation of the BERT Syntax



Experimental Results [Marco-F1]

(Wu et al., 2020)

# Towards Rationale

> **Question:** who wrote the film howl's moving castle?
>
> **Passage:** Howl's Moving Castle is a 2004 Japanese animated fantasy film written and directed by Hayao Miyazaki. It is based on the novel of the same name, which was written by Diana Wynne Jones. The film was produced by Toshio Suzuki.
>
> **Answer:** Hayao Miyazaki
>
> _____
>
> **(1) Sentence Selection**
> Howl's Moving Castle is a 2004 Japanese animated fantasy film written and directed by Hayao Miyazaki.
> **(2) Referential Equality**
> the film howl's moving castle = Howl's Moving Castle
> **(3) Entailment**
> X is a 2004 Japanese animated fantasy film written and directed by ANSWER. ⊢ ANSWER wrote X.

(Lamm et al., 2020)