



# Big Model and Brain Science

Ming Li, Zhengyan Zhang,  
Yusheng Su, Yujia Qin

THUNLP



# The Magic of Language Shared by Brain and PLM

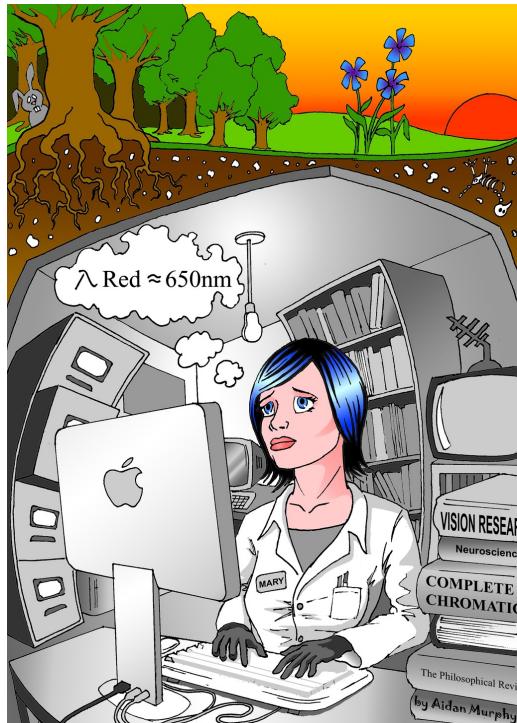
THUNLP



# Mary's room: A philosophical thought experiment

**Mary** is a scientist who exists in a black and white world where she

- has extensive access to physical descriptions of color,
- but no actual human perceptual experience of color.



One day...

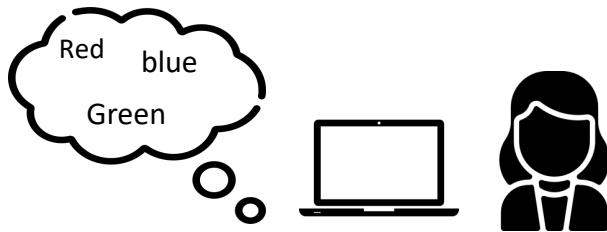
Mary walks out and she sees a blue sky!

At that moment she learns something that all her studies couldn't tell her.  
She learns what it feels to see color.

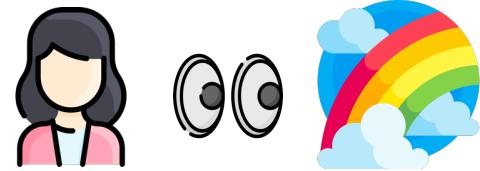


# What does it means, to know?

“Black and white” Mary



vs.



Symbolic/Language-derived representation

Pink note delicious

Artificial intelligence

Embodied/Sensory-derived representation



vs.

Human

***The computer is Mary in the black and white room.***

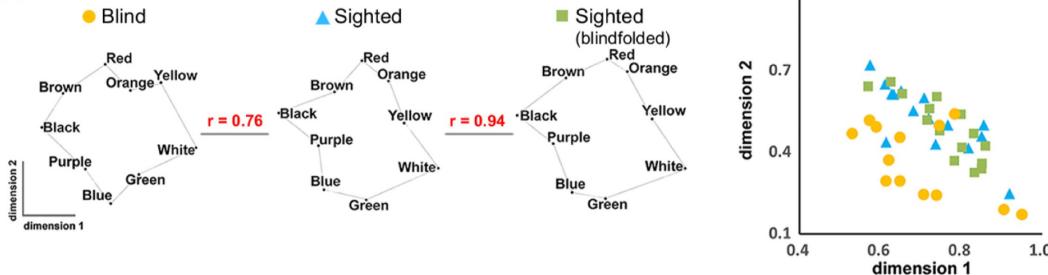
***The human is when she walks out.***



# Color knowledge for the blind and the sighted

## A. Color terms

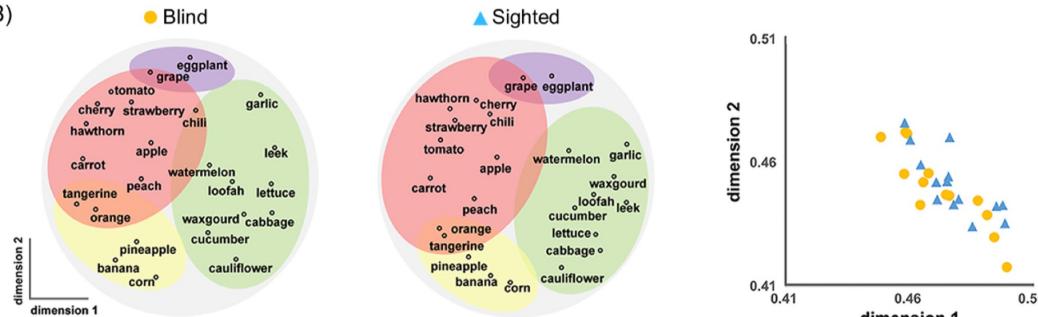
(A)



- **High degree of similarity** between blind and sighted
- Greater degree of individual variation in blind

## B. Object color properties

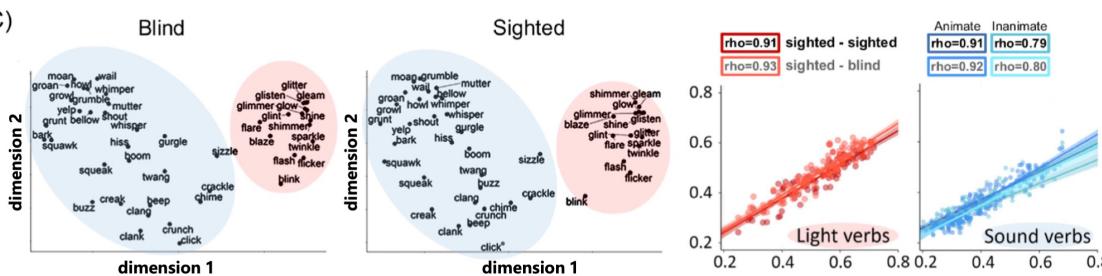
(B)



Humans can form knowledge space about a particular sensory property in the complete absence of that sensory experience

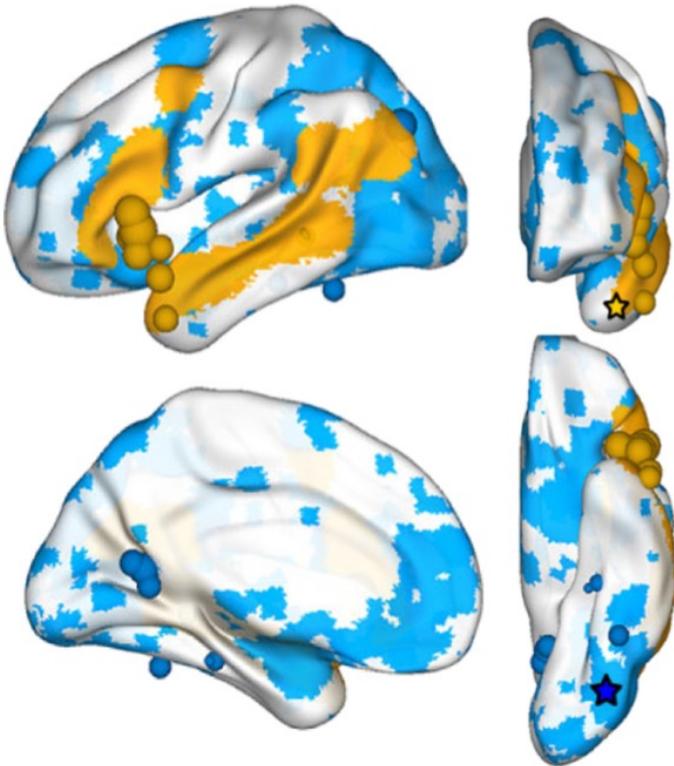
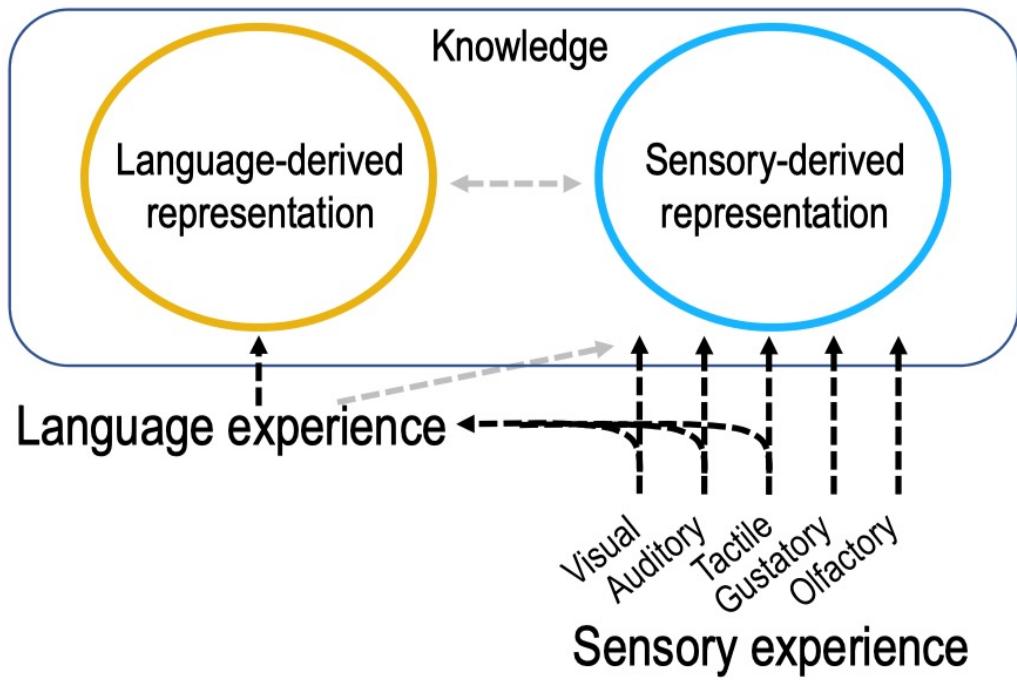
## C. Other visual-specific properties

(C)





# How does human brain code knowledge?





# Shared computational principles for language processing

## Principle 1:

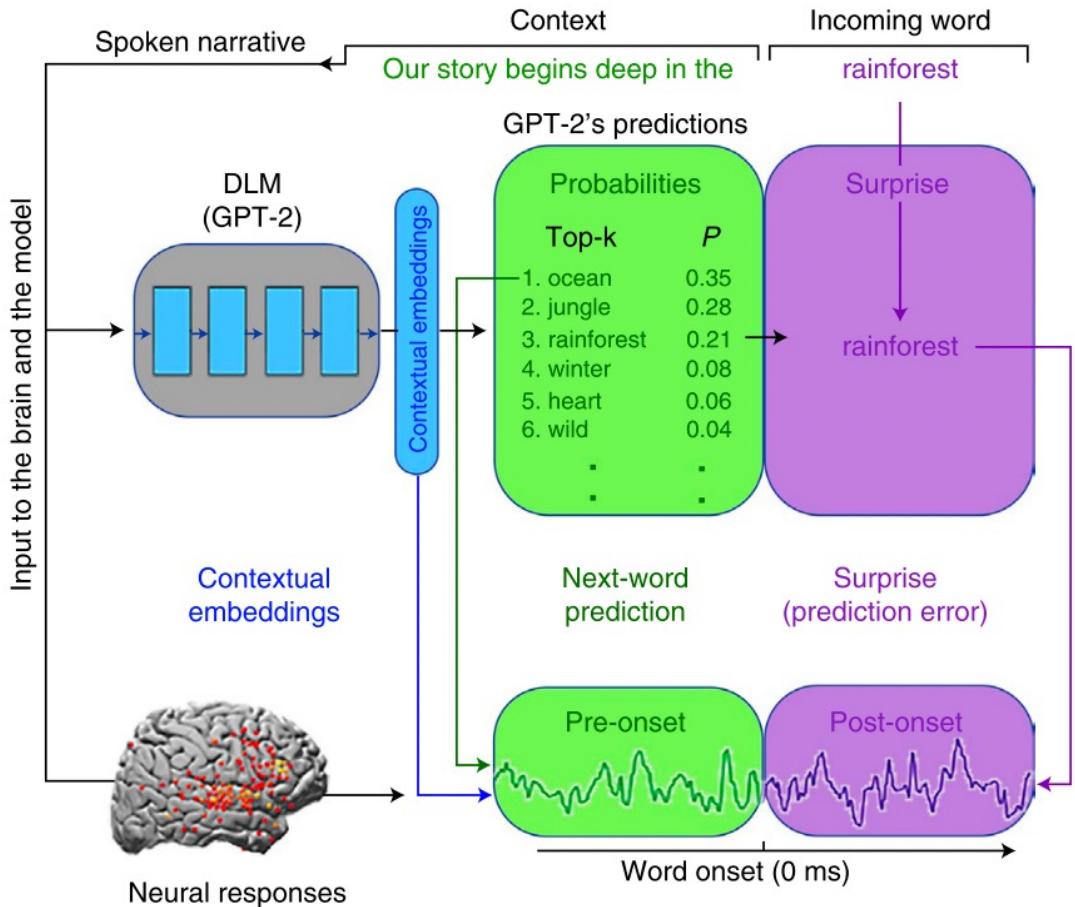
next-word prediction before word onset.

## Principle 2:

pre-onset predictions are used to calculate post-word-onset surprise.

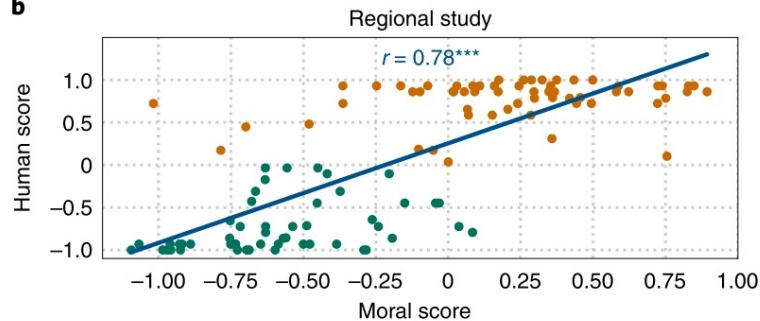
## Principle 3:

contextual vectorial representation in the brain.

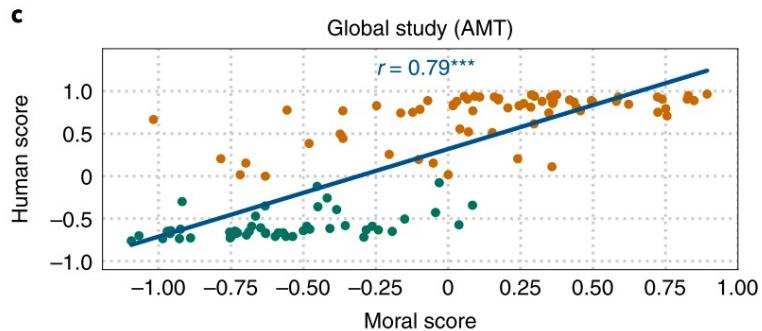


# Revealing the magic of language – Function

b

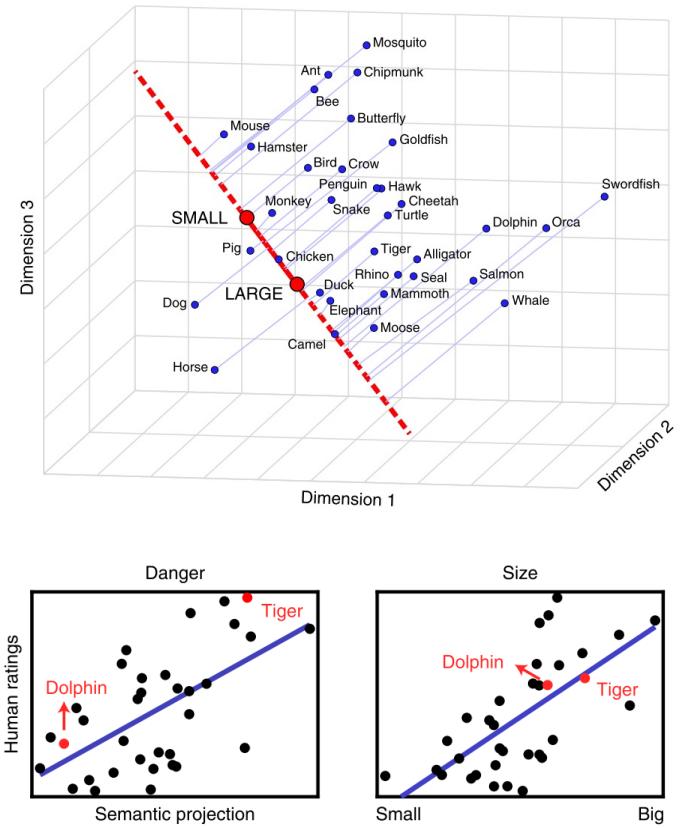
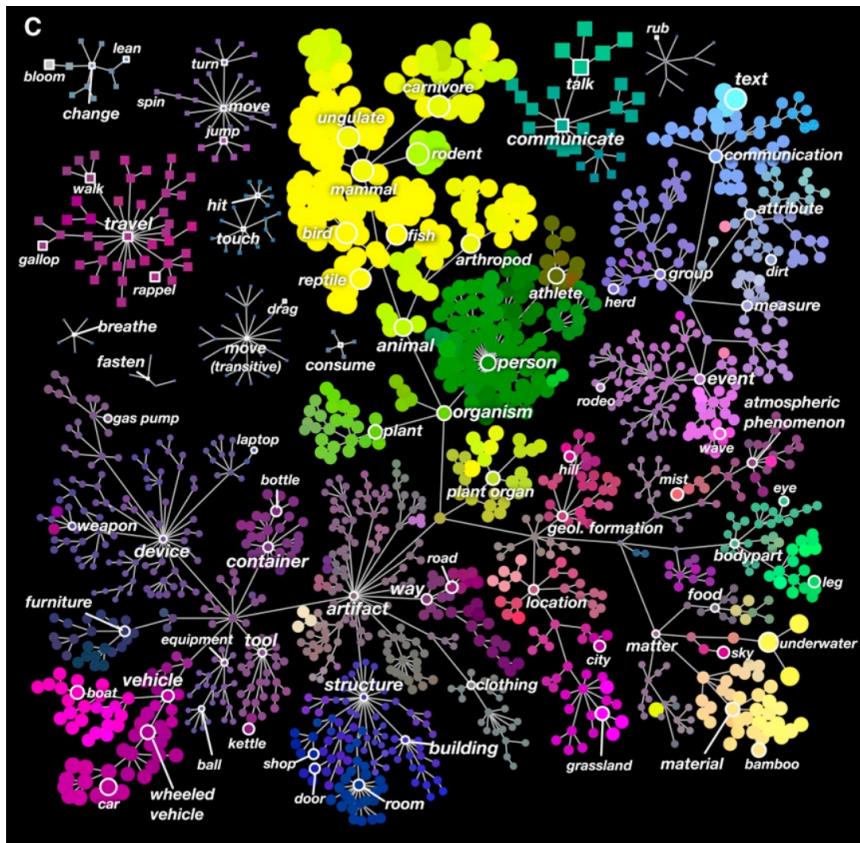


c





# Revealing the magic of language – Representation

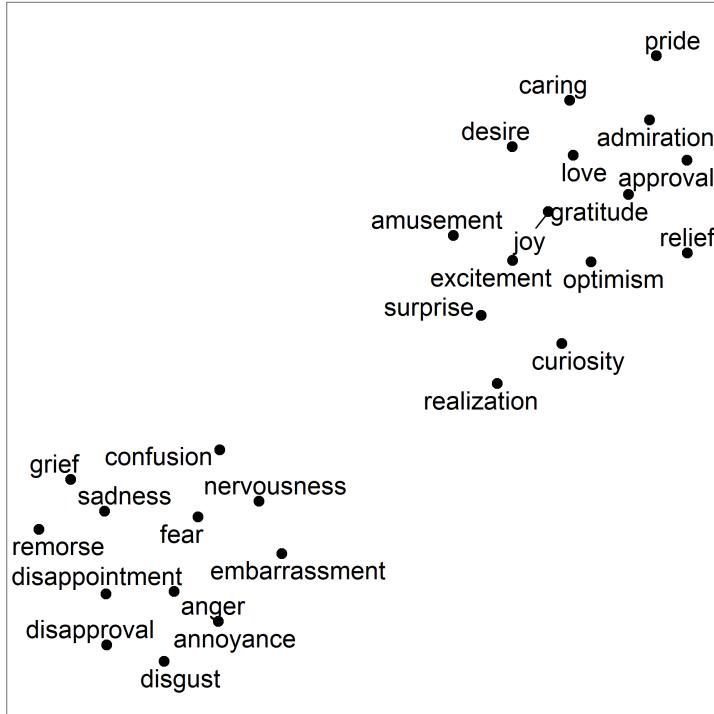


Note that the semantic representations derived from language input **do not possess** feelings or experiences of the world. Such representations **do reflect** perceptual (size), abstract (danger), and even affective (arousal and valence) properties of concepts

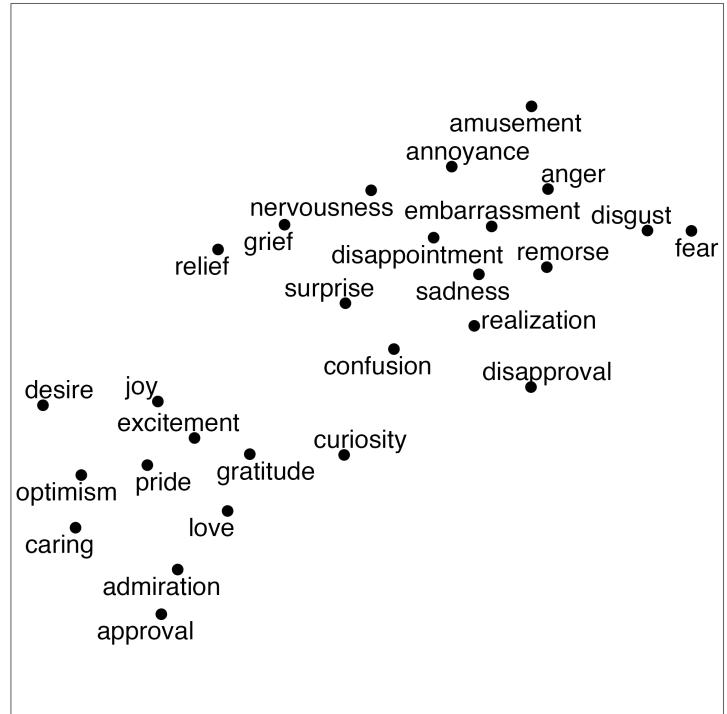


# Revealing the magic of language – Representation

## Mental Representation



## Semantic Representation



Note that the semantic representations derived from language input **do not possess** feelings or experiences of the world. Such representations **do reflect** perceptual (size), abstract (danger), and even affective (arousal and valence) properties of concepts



# Revealing the magic of language – Structure

## Language Stimuli

Pereira2018

"Beekeeping encourages the conservation of local habitats. It is in every beekeeper's interest..."

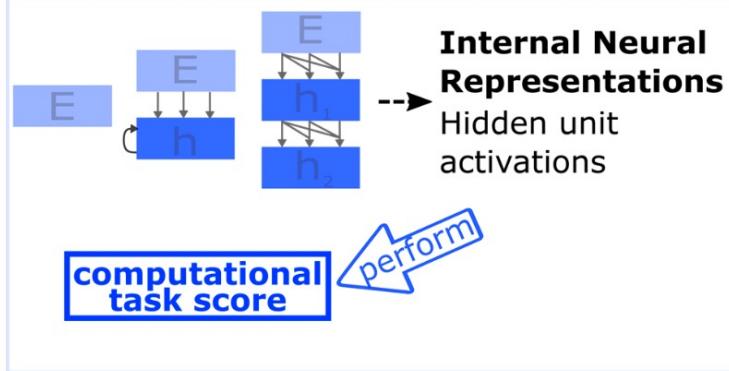
Fedorenko2016

"Alex was tired so he took a nap."

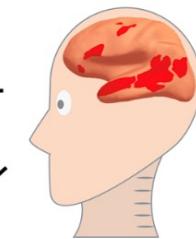
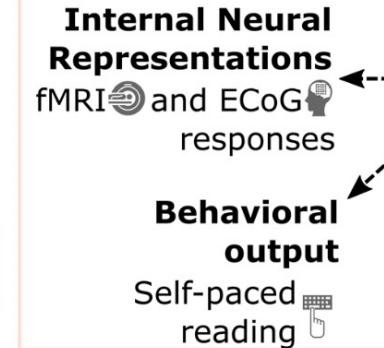
Blank2014

"If you were to journey to the North of England, you would come to a valley that is surrounded by moors as high as mountains. It is in this valley where you..."

## Models

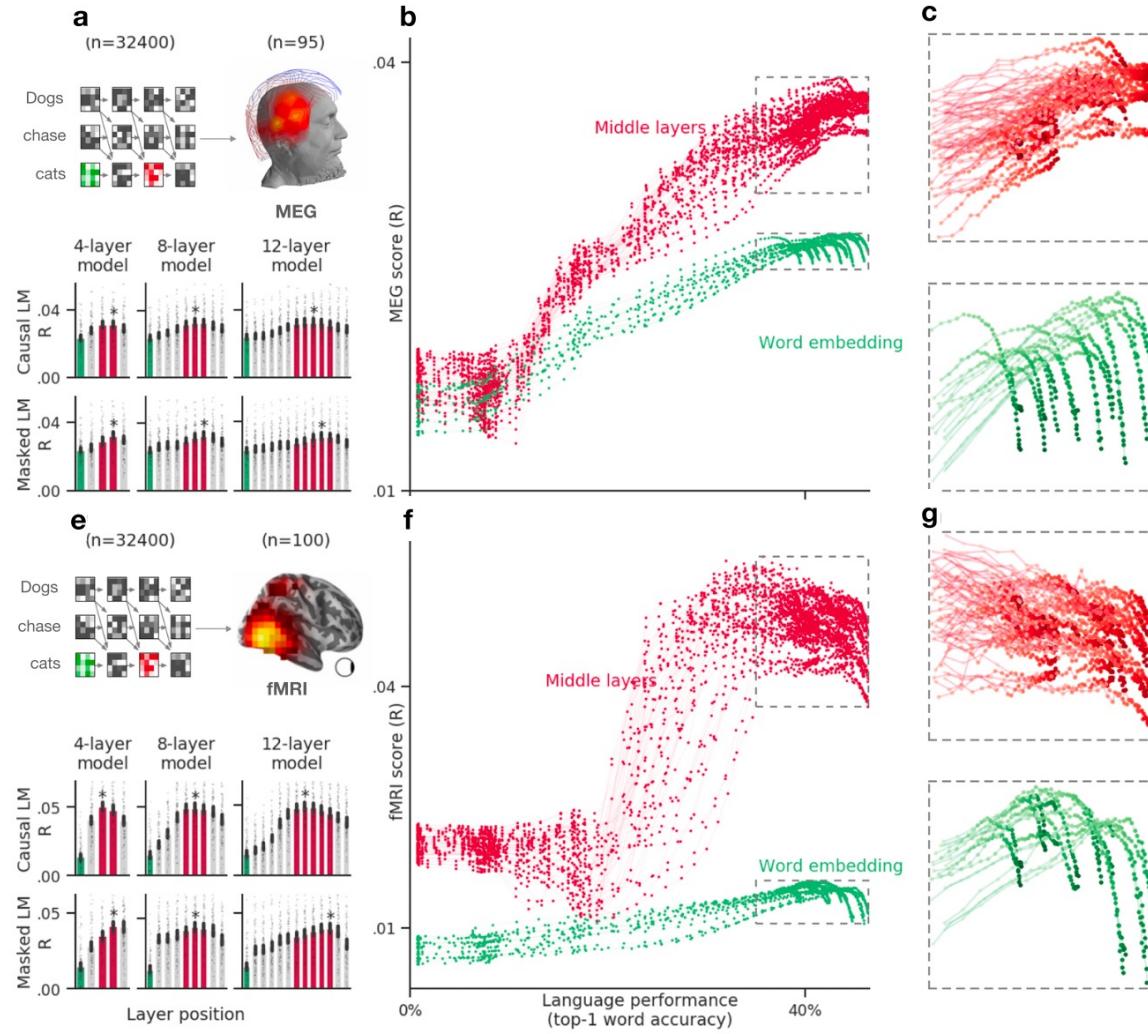


## Humans



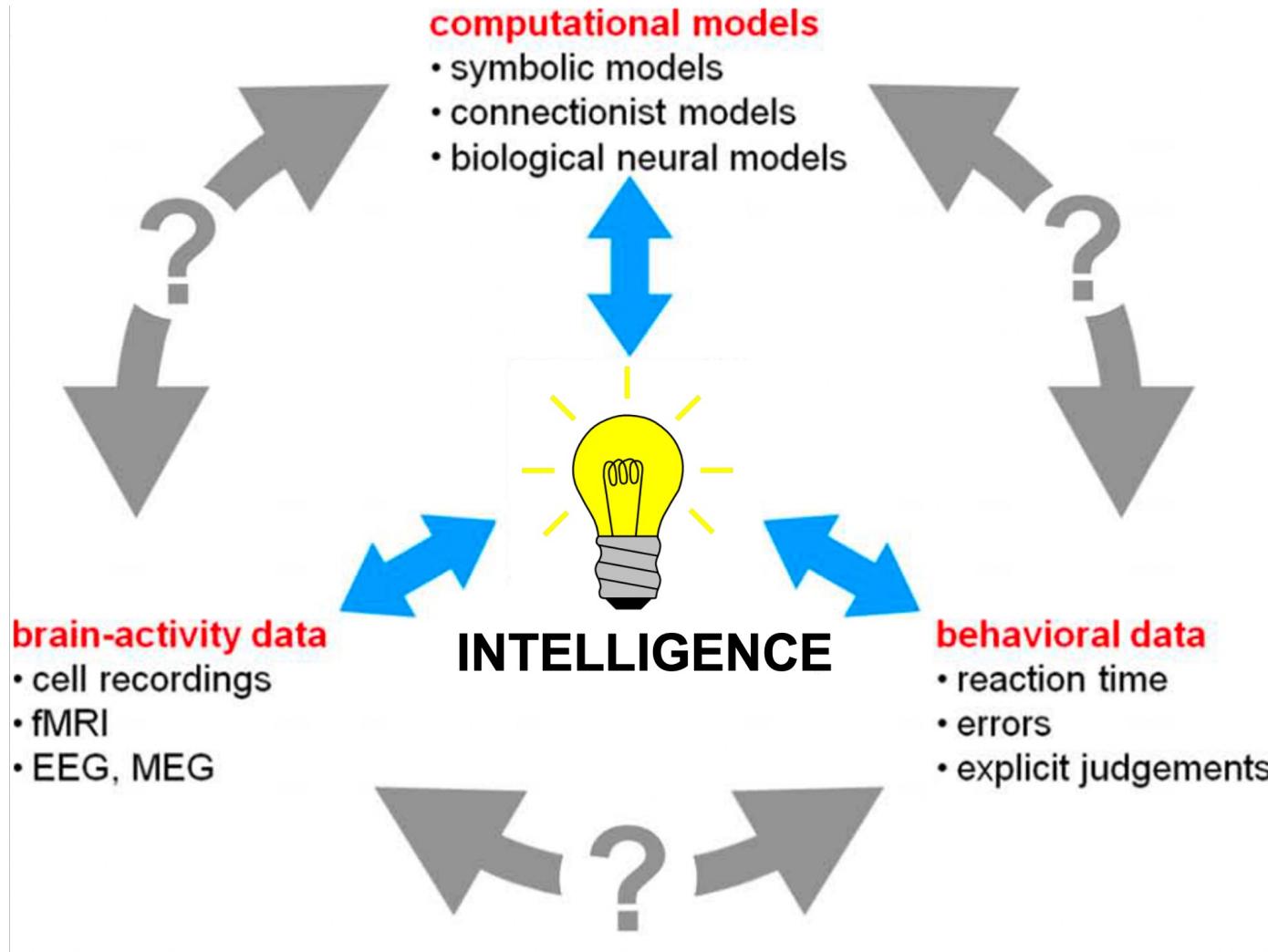


# Revealing the magic of language – Structure





# The next question – Towards an understanding of intelligence



\* Figure is adjusted from –

Kriegeskorte, N., Mur, M. & Bandettini, P. Representational similarity analysis - connecting the branches of systems neuroscience. *Front. Syst. Neurosci.* **2**, 1–28 (2008). 13



# Neuron Activation

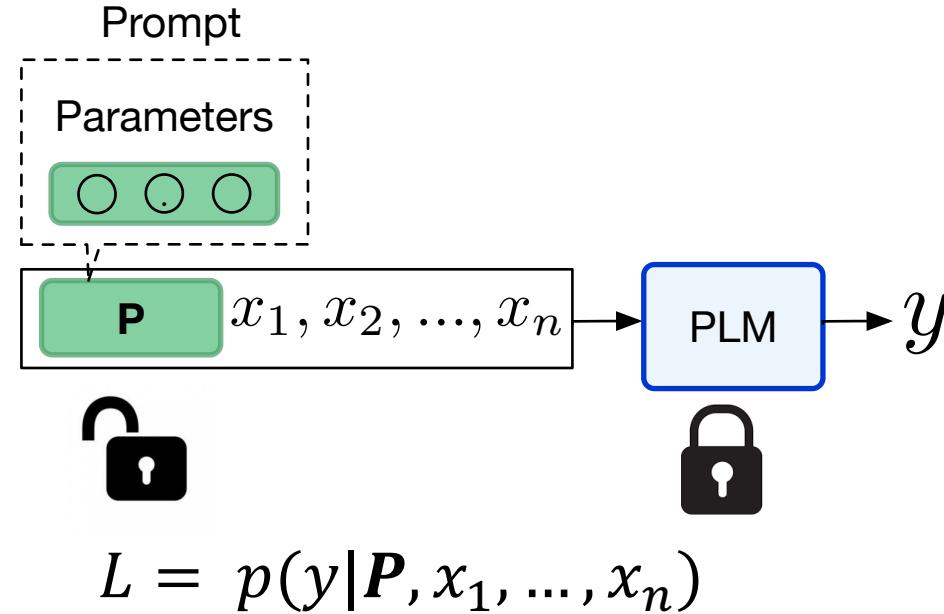
THUNLP

# Transferability Indicator



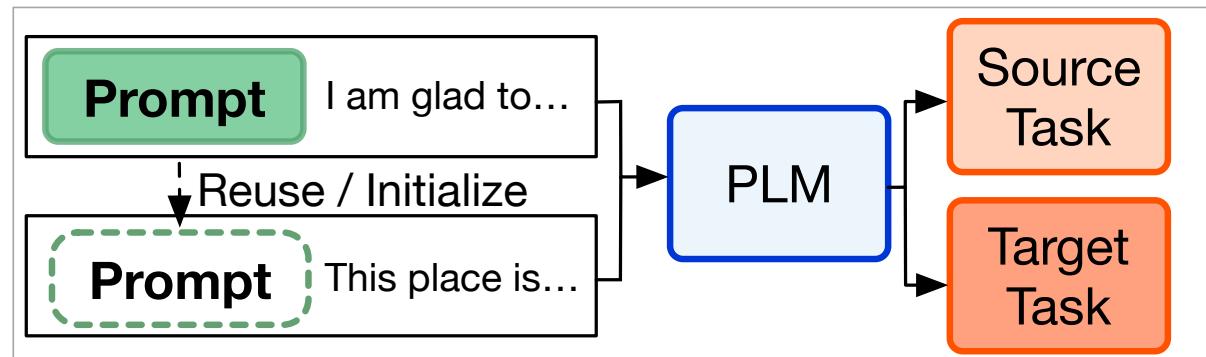
# Recap prompt tuning

- Training:



- Transferability:

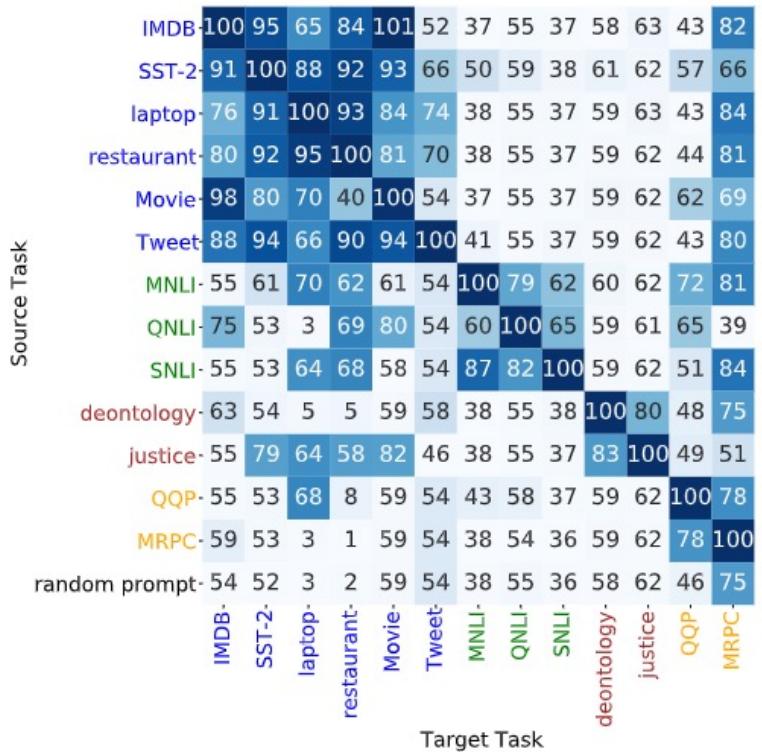
## Cross-Task Transfer





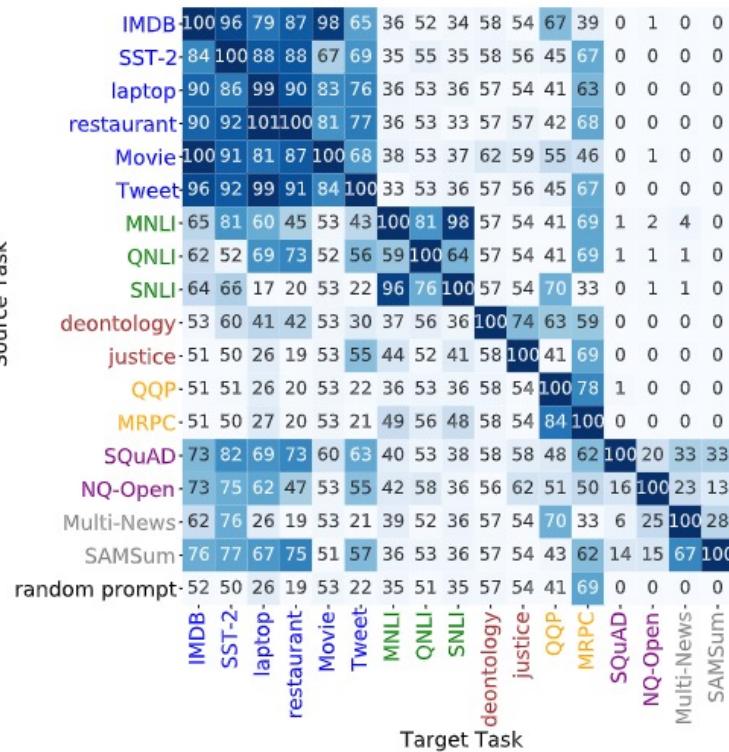
# Prompt transfer

- Cross-Task Transfer (Zero-shot)
  - For the tasks *within the same type*, transferring prompts between them can generally perform well.



(a) RoBERTa<sub>LARGE</sub>

(Relative Performance)

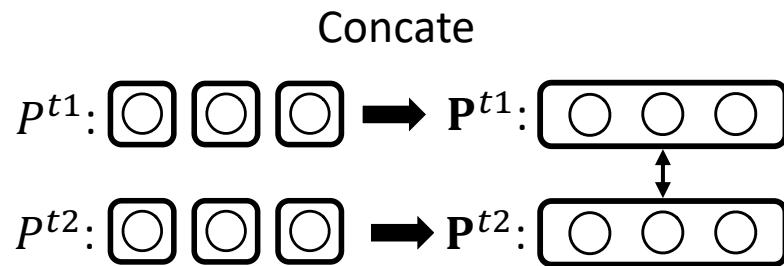


(b) T5<sub>XXL</sub>



# Transferability indicator

- Motivation
  - Explore why the soft prompts can transfer across tasks and what decides the transferability between them
- Embedding Similarity
  - Euclidean similarity
  - Cosine similarity



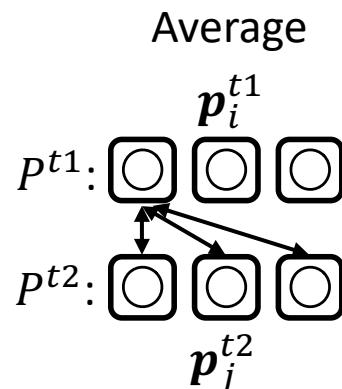
$$E_{\text{concat}}(P^{t1}, P^{t2}) = \frac{1}{1 + \|\mathbf{P}^{t1} - \mathbf{P}^{t2}\|}$$

$$C_{\text{concat}}(P^{t1}, P^{t2}) = \frac{\mathbf{P}^{t1} \cdot \mathbf{P}^{t2}}{\|\mathbf{P}^{t1}\| \|\mathbf{P}^{t2}\|}$$



# Transferability indicator

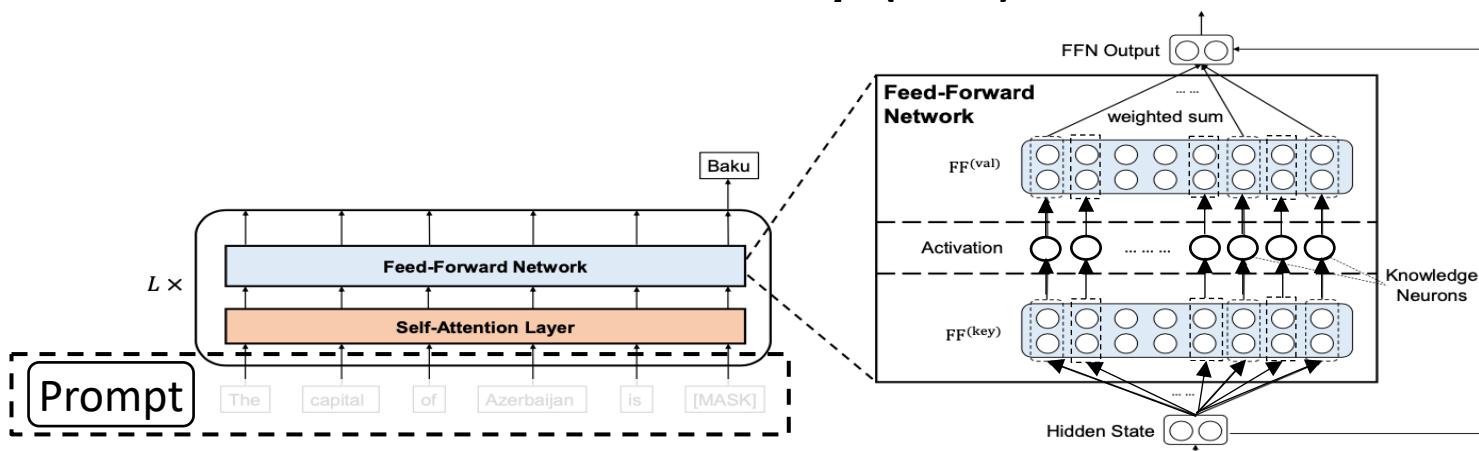
- Motivation
  - Explore why the soft prompts can transfer across tasks and what decides the transferability between them
- Embedding Similarity
  - Euclidean similarity
  - Cosine similarity



$$E_{\text{average}}(P^{t1}, P^{t2}) = \frac{1}{1 + \frac{1}{l^2} \sum_{i=1}^l \sum_{j=1}^l \|\mathbf{p}_i^{t1} - \mathbf{p}_j^{t2}\|}$$
$$C_{\text{average}}(P^{t1}, P^{t2}) = \frac{1}{l^2} \sum_{i=1}^l \sum_{j=1}^l \frac{\mathbf{p}_i^{t1} \cdot \mathbf{p}_j^{t2}}{\|\mathbf{p}_i^{t1}\| \|\mathbf{p}_j^{t2}\|}$$

# Transferability indicator

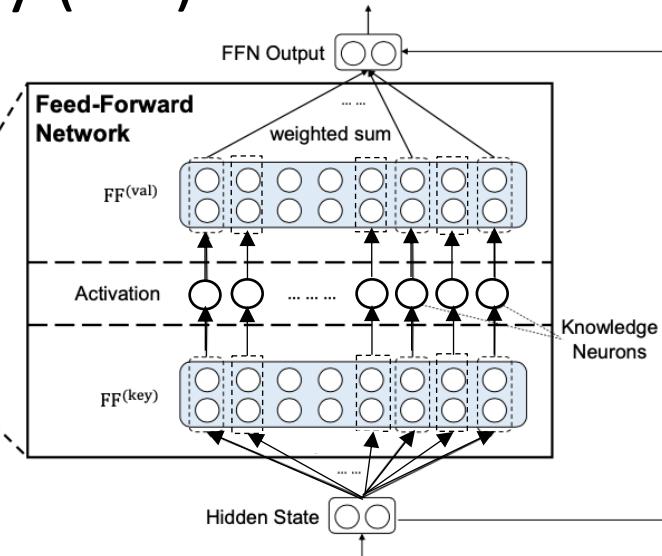
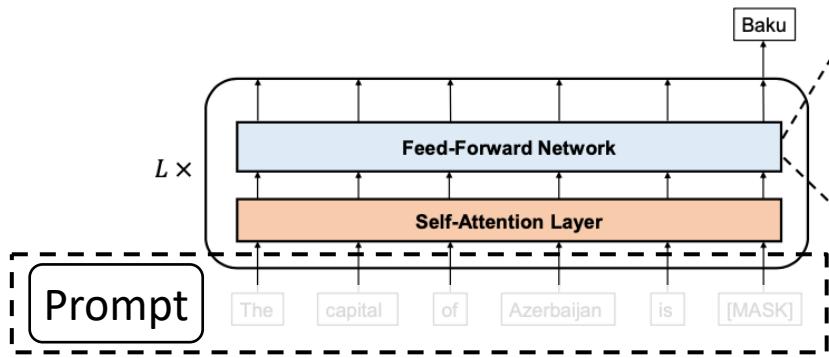
- Motivation
  - Explore why the soft prompts can transfer across tasks and what decides the transferability between them
- Embedding Similarity
  - Euclidean similarity
  - Cosine similarity
- Model Stimulation Similarity (ON)





# Transferability indicator

- Model Stimulation Similarity (ON)
  - Activated Neurons:

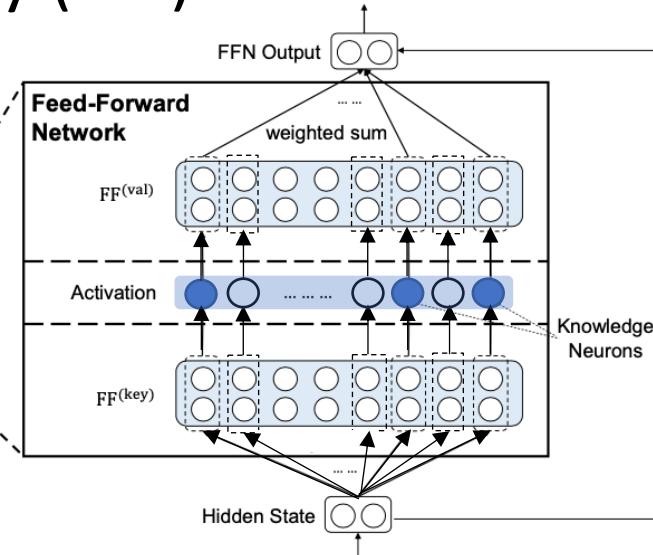
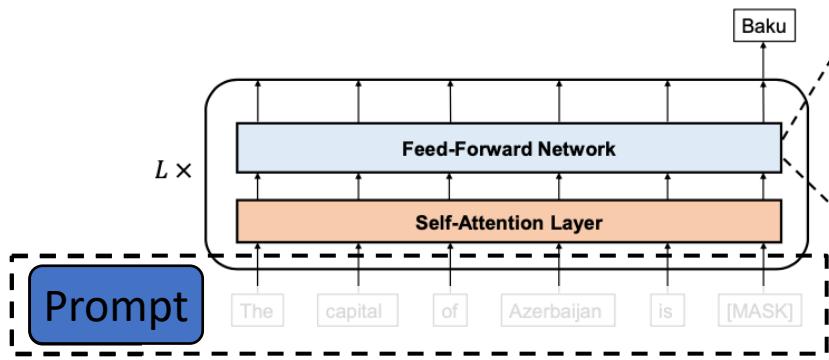


$$\text{FFN}(\mathbf{x}) = \max(\mathbf{x}W_1^\top + \mathbf{b}_1, 0)W_2 + \mathbf{b}_2,$$
$$\text{AS}(P) = [\mathbf{s}_1; \mathbf{s}_2; \dots; \mathbf{s}_L]$$



# Transferability indicator

- Model Stimulation Similarity (ON)
  - Activated Neurons:



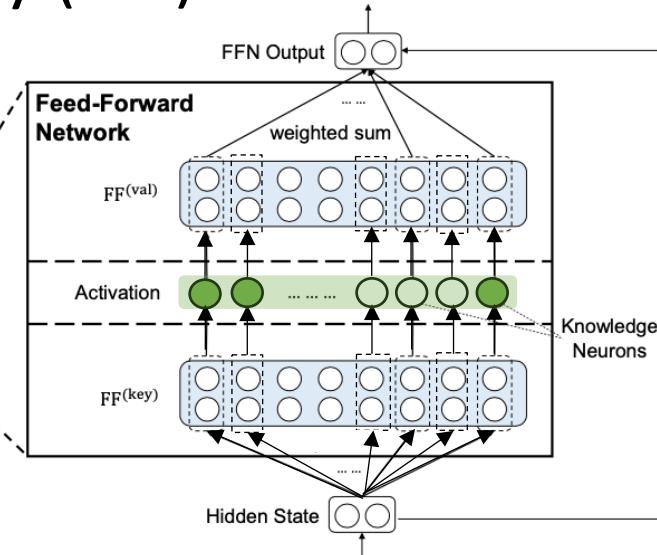
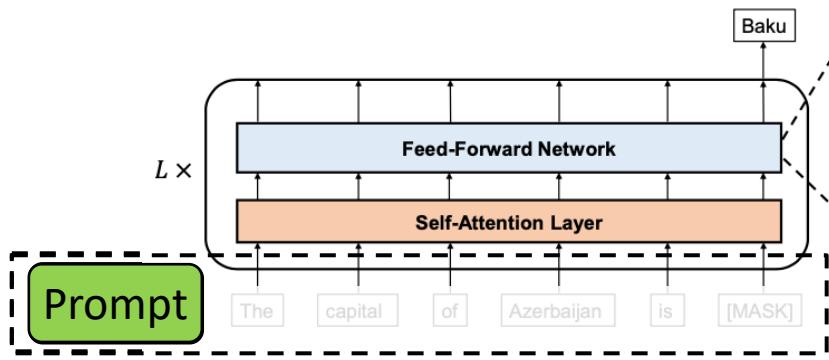
$$\text{FFN}(\mathbf{x}) = \max(\mathbf{x}W_1^\top + \mathbf{b}_1, 0)W_2 + \mathbf{b}_2,$$
$$\text{AS}(P) = [\mathbf{s}_1; \mathbf{s}_2; \dots; \mathbf{s}_L]$$





# Transferability indicator

- Model Stimulation Similarity (ON)
  - Activated Neurons:



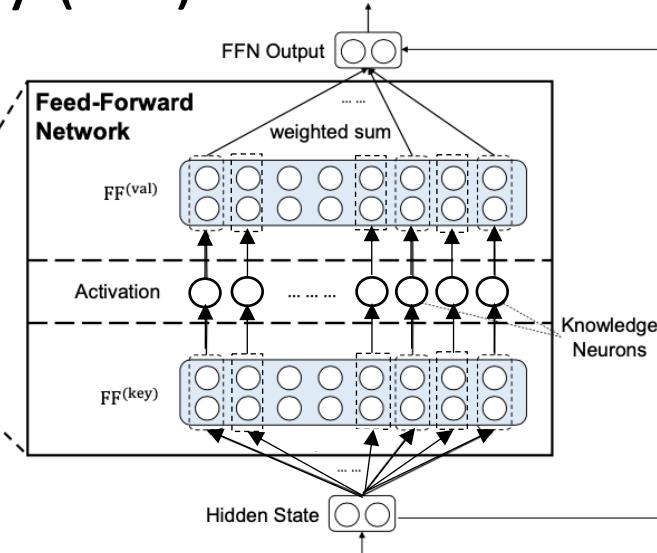
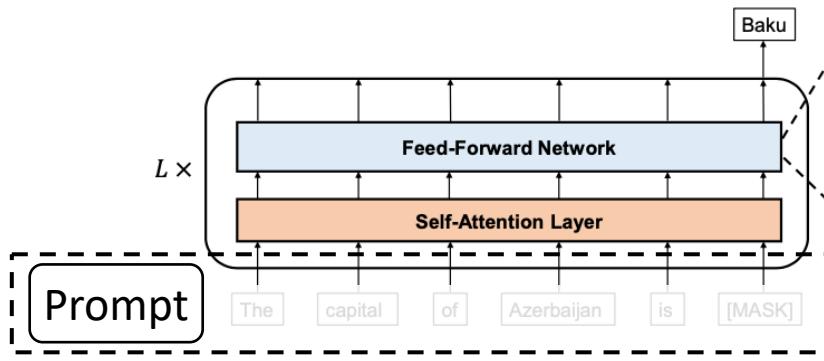
$$\text{FFN}(\mathbf{x}) = \max(\mathbf{x}W_1^\top + \mathbf{b}_1, 0)W_2 + \mathbf{b}_2,$$
$$\text{AS}(P) = [\mathbf{s}_1; \mathbf{s}_2; \dots; \mathbf{s}_L]$$





# Transferability indicator

- Model Stimulation Similarity (ON)
  - Activated Neurons:



$$\text{FFN}(\mathbf{x}) = \max(\mathbf{x}W_1^\top + \mathbf{b}_1, 0)W_2 + \mathbf{b}_2,$$
$$\text{AS}(P) = [\mathbf{s}_1; \mathbf{s}_2; \dots; \mathbf{s}_L]$$

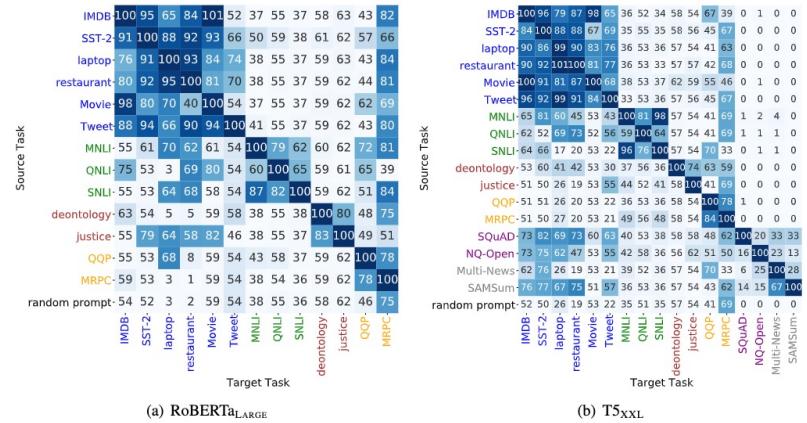


$$\text{ON}(P^{t_1}, P^{t_2}) = \frac{\text{AS}(P^{t_1}) \cdot \text{AS}(P^{t_2})}{\|\text{AS}(P^{t_1})\| \|\text{AS}(P^{t_2})\|}$$



# Transferability indicator

- Similarity
  - Embedding Similarity
    - Euclidean similarity:
    - Cosine similarity:
  - Model Stimulation Similarity (ON)
- Zero-shot Task Transferability



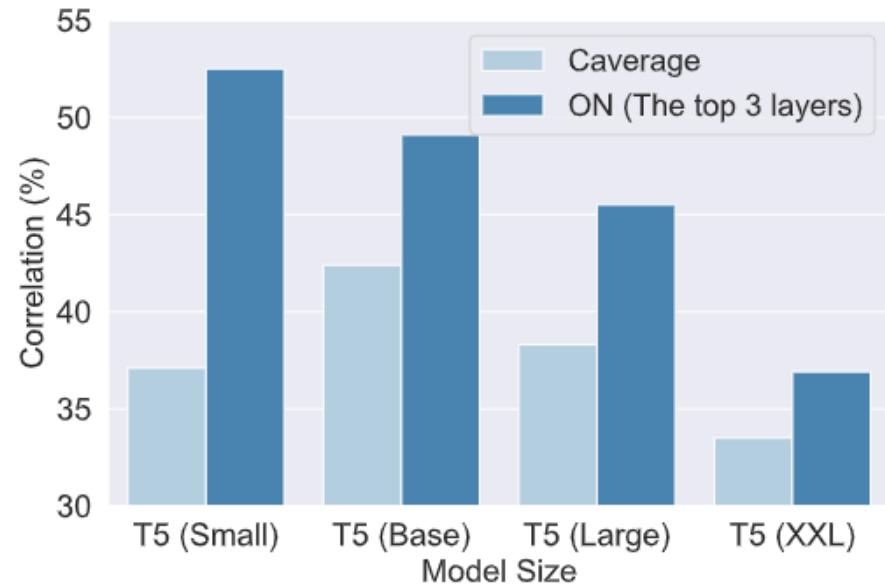
← Spearman's correlation →



# Transferability indicator

- Model Stimulation Similarity (ON)
  - ON has the higher Spearman's correlation with the transferability
  - ON works worse on the larger PLMs because of the higher redundancy [1]

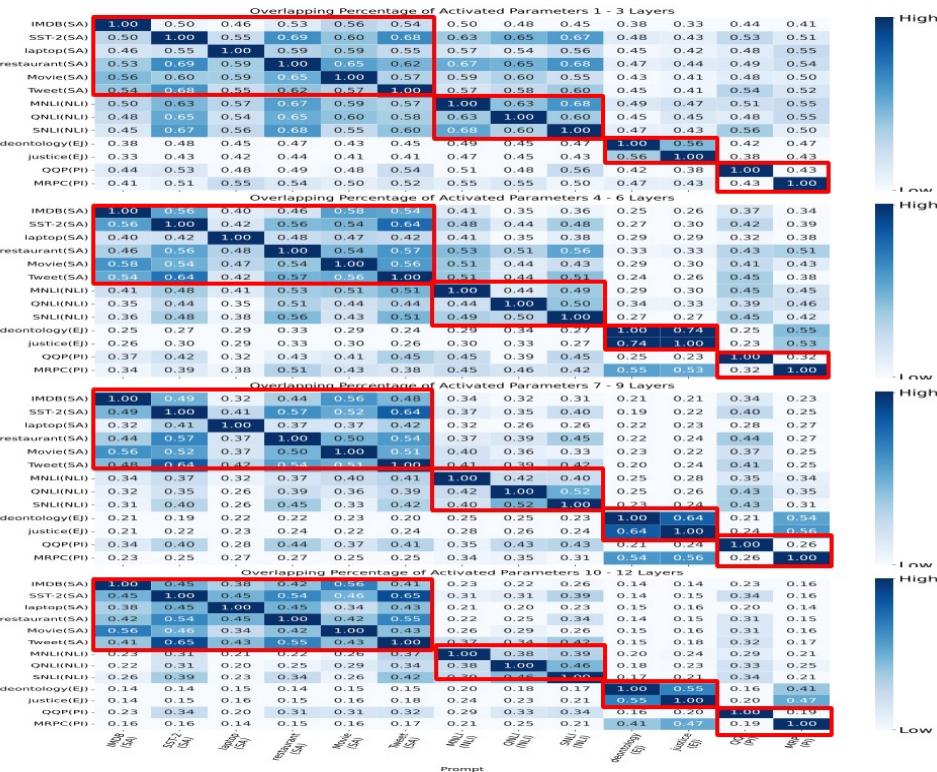
Metric	Model	
	RoBERTa <sub>LARGE</sub>	T5 <sub>XXL</sub>
E <sub>concat</sub>	22.6	12.9
E <sub>average</sub>	2.8	-2.5
C <sub>concat</sub>	24.8	31.6
C <sub>average</sub>	44.7	33.5
ON	<b>49.7</b>	<b>36.9</b>





# Activated neurons in a PLM

- Distribution of Activated Neuron
  - The activated neurons are common in the bottom layers but more task-specific in top layers.



1 – 3 layers

4 – 6 layers

7 – 9 layers

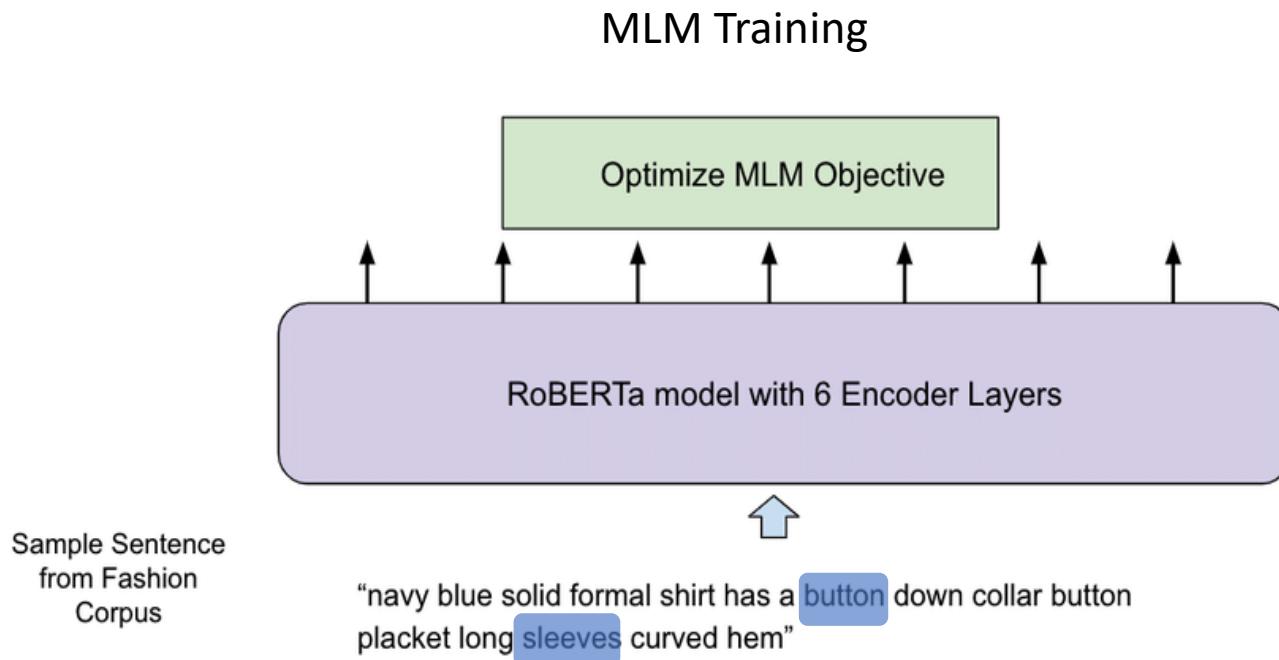
10 – 12 layers

Activated Neurons Can  
Reflect Human-Like  
Emotional Attributes



# Question

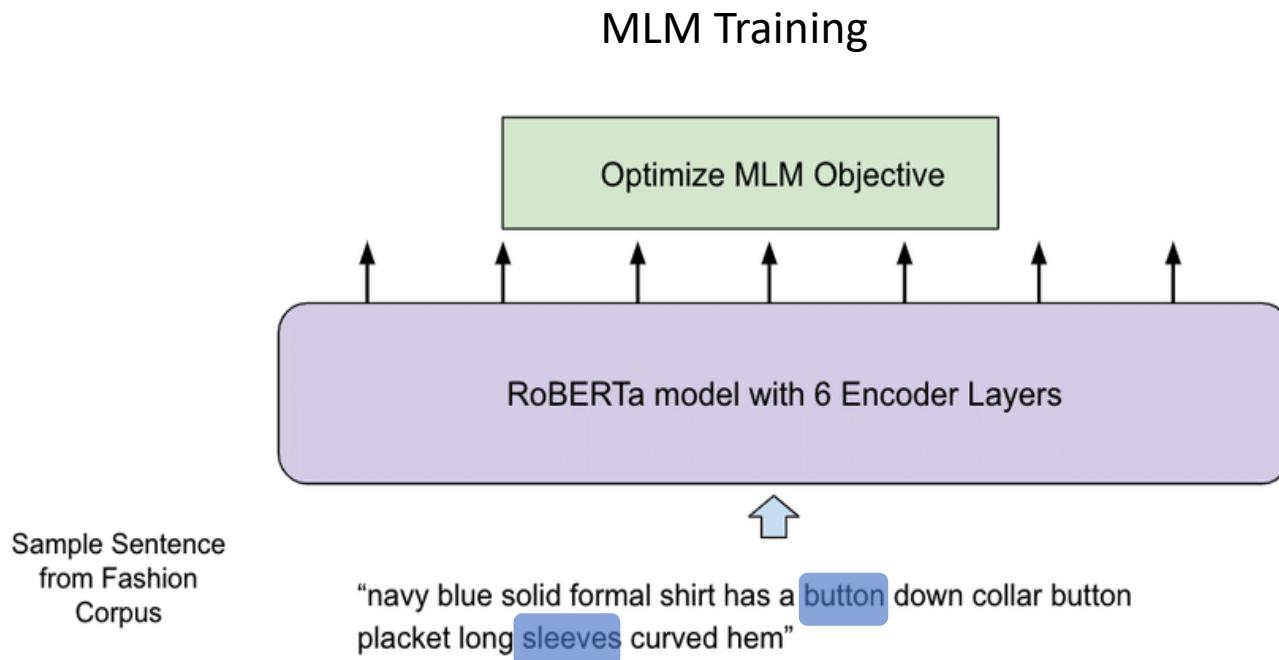
- Whether PLMs can learn human-like emotional attributes during the pre-training ?





# Question

- Whether PLMs can learn human-like emotional attributes during the pre-training ?



Can learn human-like emotional attributes ?



# How do humans recognize different emotions ?

- Human

## Rating of Emotional Attributes

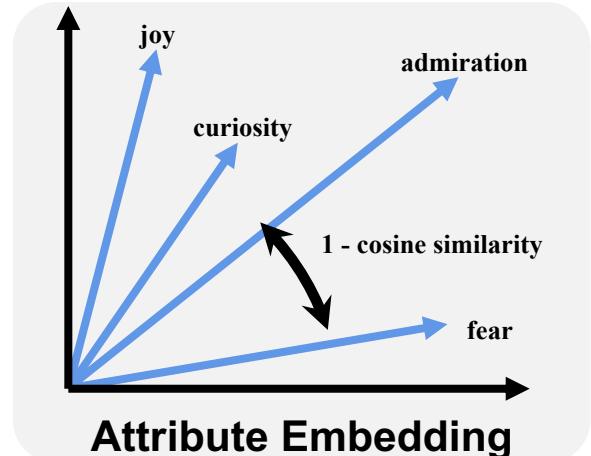
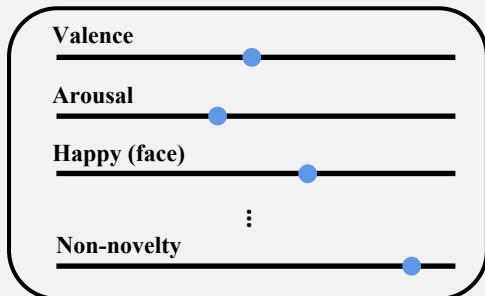
### Attributes:

- [Arousal, Valence, Happy, Anger, Sad, Fear, Surprise, Disgust, Control, Fairness, Self-related, Other-related, Expectedness, Non-novelty]: 14

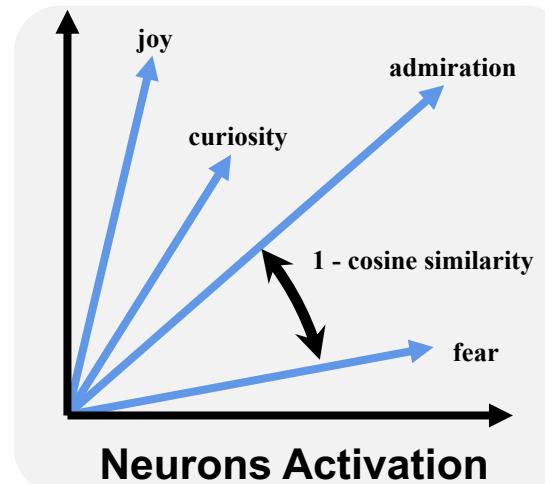
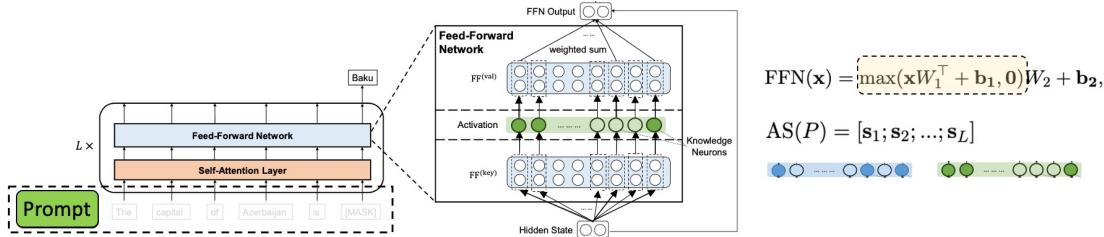


N = 30 / 30 / 299

admiratio



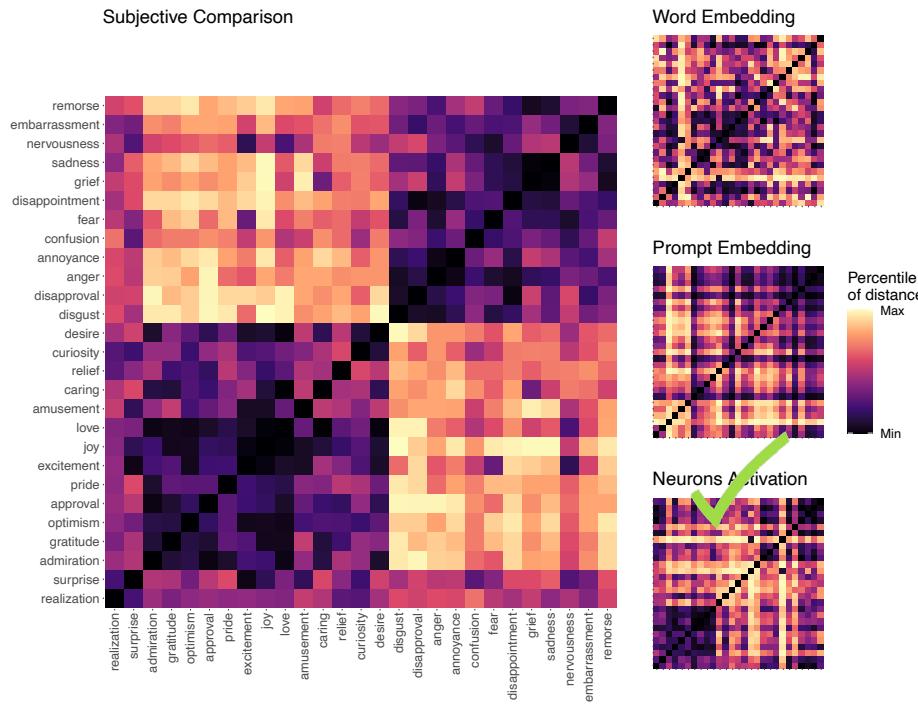
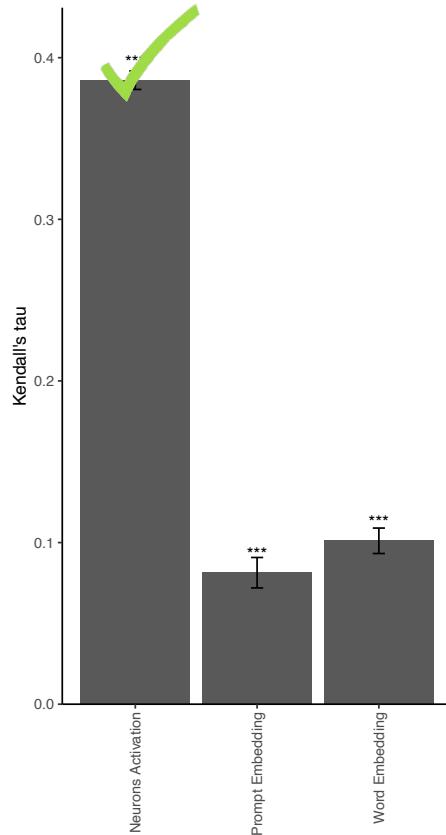
- PLM (Activated Neurons)





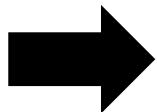
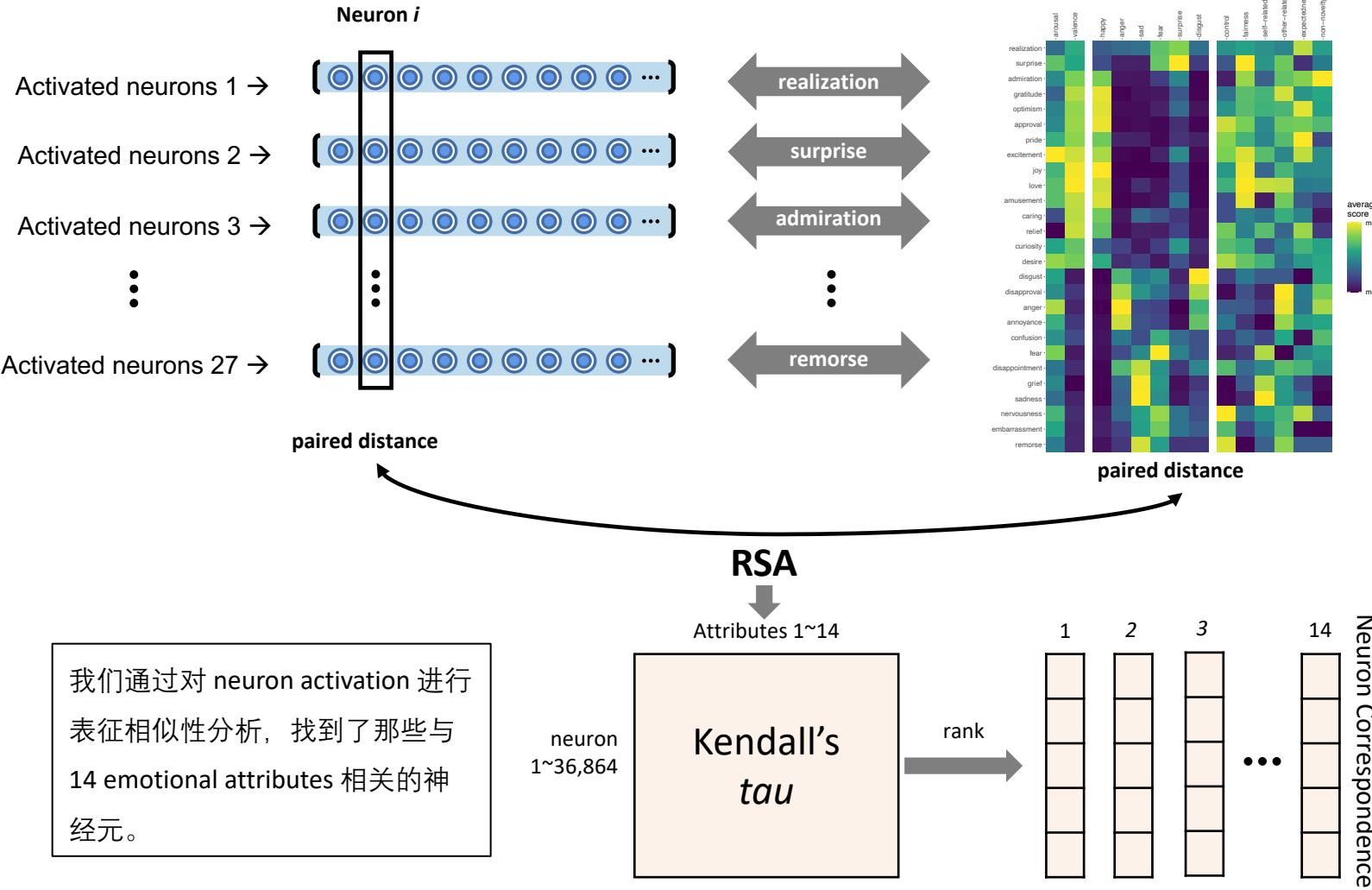
# Correlation

- Correlation:
  - Represent 27 emotions with human attributes and activated neurons





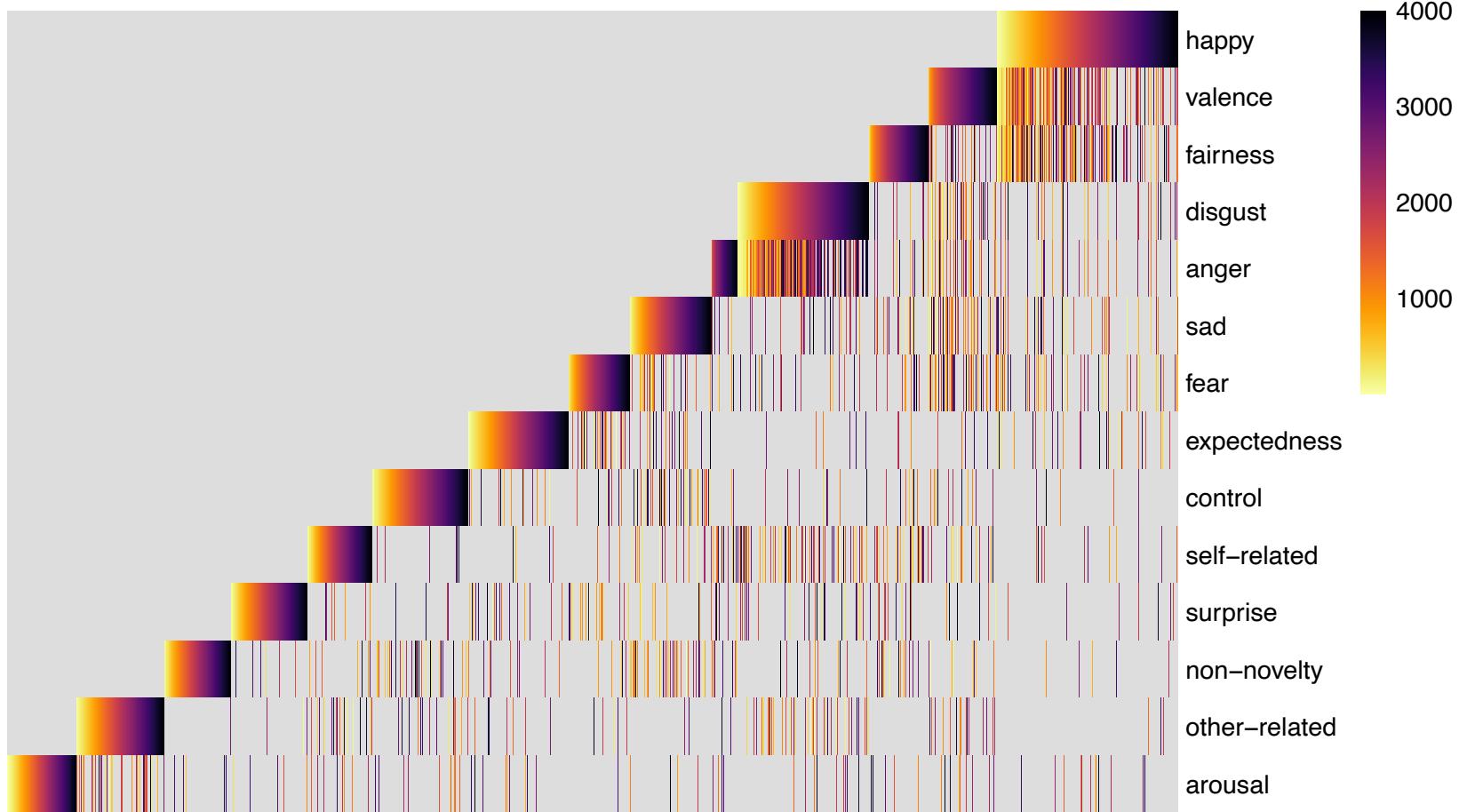
# Activated neurons for every attribute



**Attribute: Activated Neurons (○○○○○○○○○○...)**

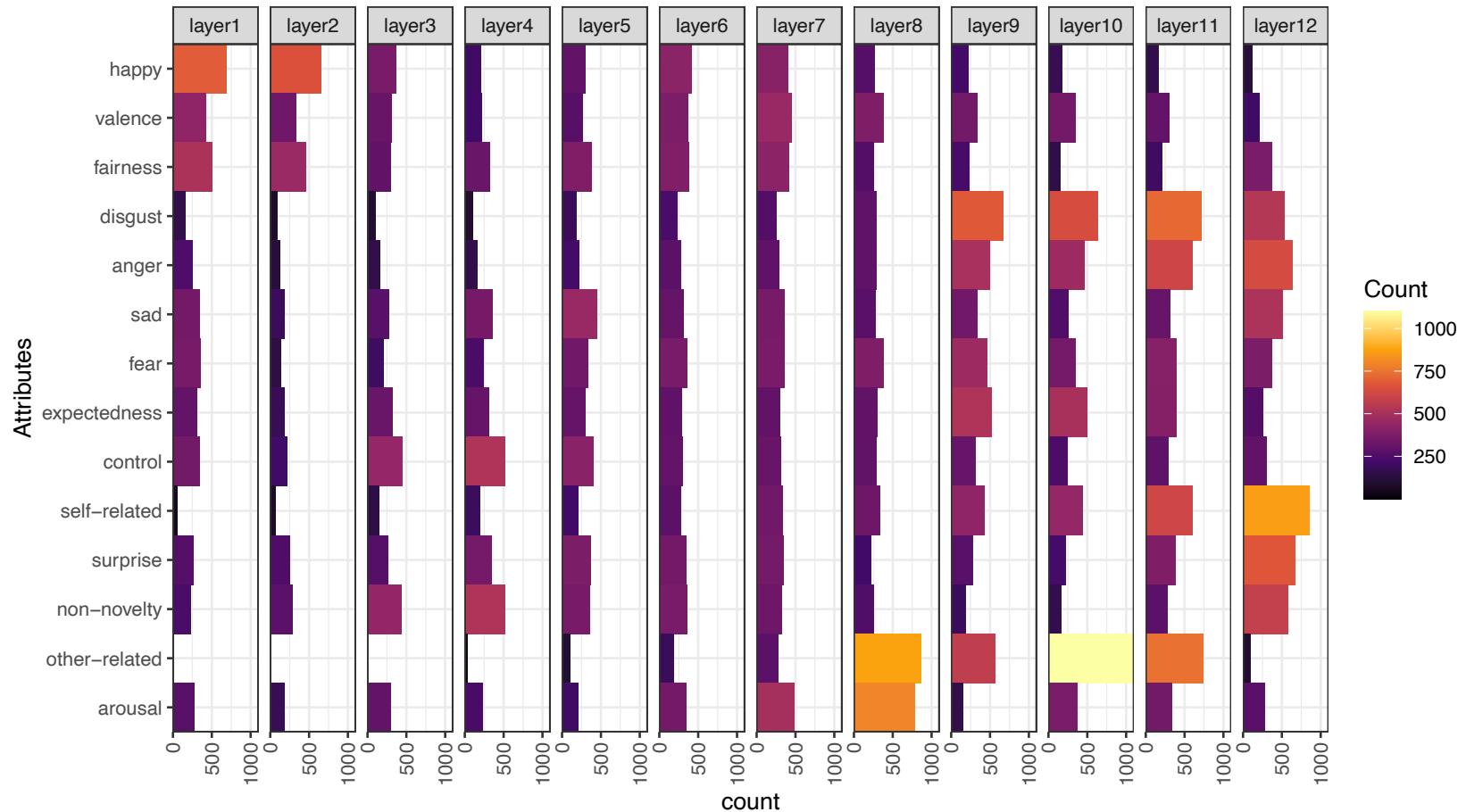


# Activated neurons for every attribute



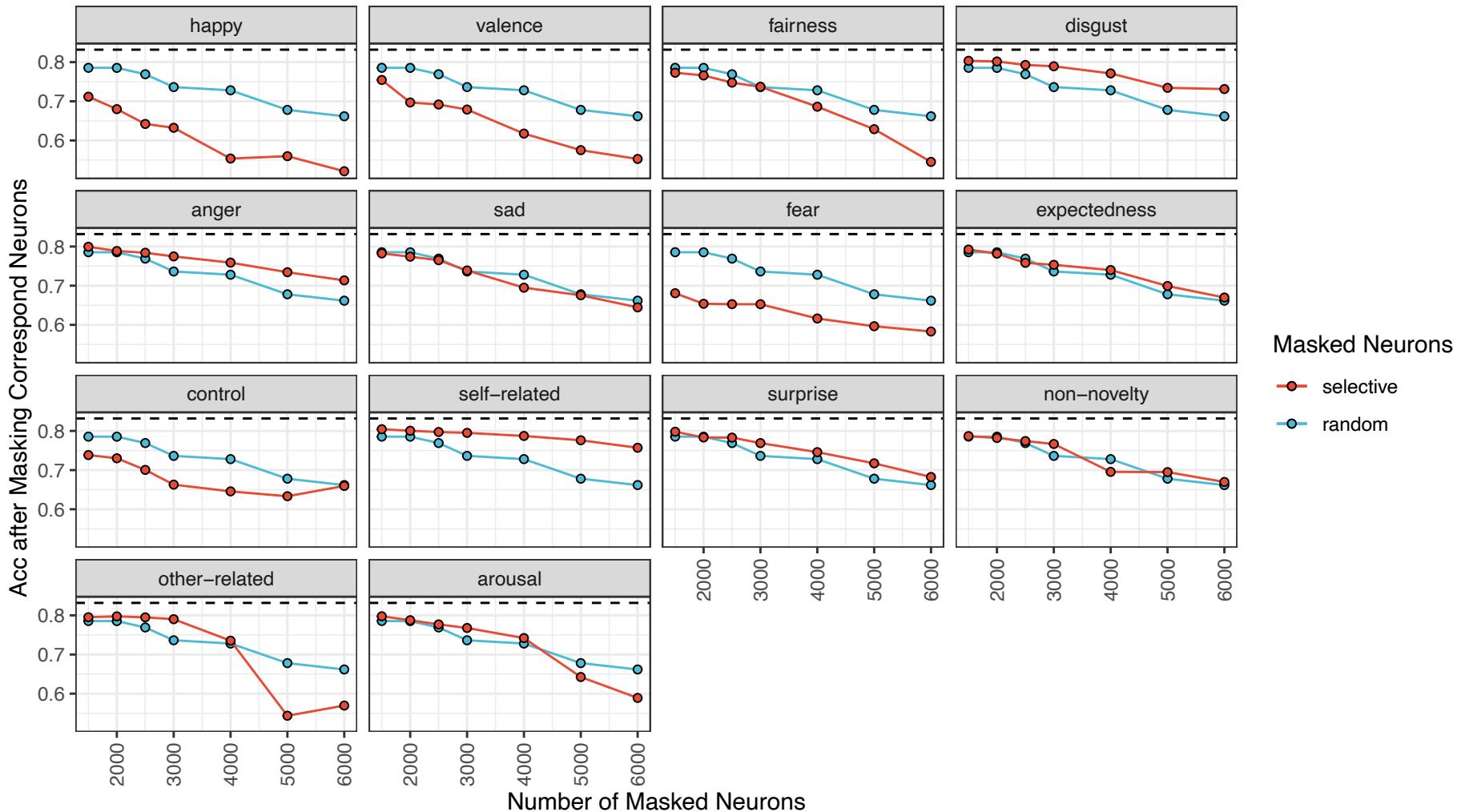


# Activated neurons for every attribute





# Remove neurons for an attribute



Masked Neurons

- selective
- random

# Demo

**thunlp/Prompt-Transferability**

On Transferability of Prompt Tuning for Natural Language Processing

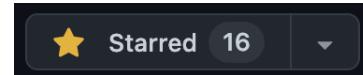
自然语言处理与  
社会人文计算实验室  
TSINGHUA UNIVERSITY

1 Contributor    0 Issues    15 Stars    0 Forks



<https://github.com/thunlp/Prompt-Transferability>

Find: Activated Neurons Demo [Colab link]



[Click]

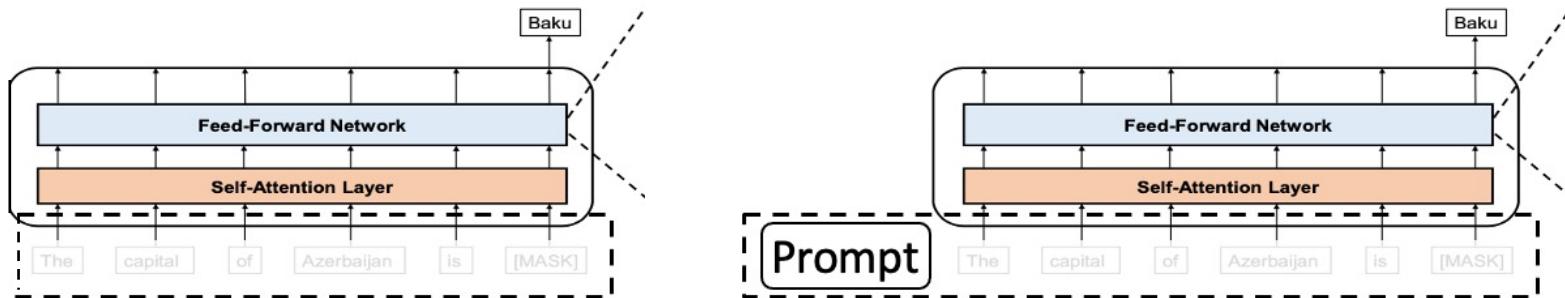


# Activated Neurons

- Load Pre-trained Language Model (Roberta)
- Load the prompts (checkpoints) - 27 Emotion Tasks

```
[4] # 27 Tasks
sentiments = ['amusement', 'excitement', 'joy', 'love', 'desire', 'optimism', 'caring',
              'pride', 'admiration', 'gratitude', 'relief', 'approval',
              'realization', 'surprise', 'curiosity', 'confusion',
              'fear', 'nervousness', 'remorse', 'embarrassment', 'disappointment',
              'sadness', 'grief', 'disgust', 'anger', 'annoyance', 'disapproval']

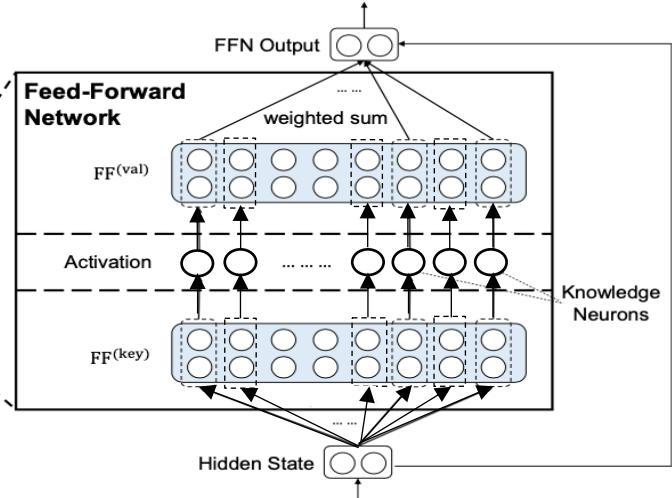
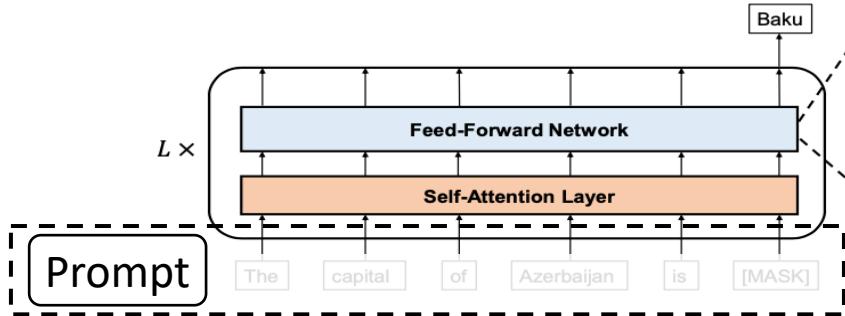
# 12 random seeds (We only use 1 seed for the demonstration)
random_seeds = [1]
# random_seeds = [1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 42, 100]
```





# Activated Neurons

- Activate Neurons



```
ckpt = torch.load(f'./checkpoint/{sentiment}-{seed}.pt')    # Load the trained prompt
prompt_emb = Parameter(ckpt['best_prompt']).to('cpu')
model.roberta.embeddings.prompt_embedding.weight.data = prompt_emb

# Forward pass to get active neuron
outputs = [[] for _ in range(nlayers)]
def save_ppt_outputs1_hook(n):
    def fn(_, __, output):
        outputs[n] = output.detach().to('cpu')      # Save the neuron of the n layer to outputs[n]
    return fn

for n in range(nlayers):
    # register_forward_hook: This hook will be used every time forward() is called
    model.roberta.encoder.layer[n].intermediate.register_forward_hook(save_ppt_outputs1_hook(n))

for sentence in loader:
    inputs = tokenizer(sentence, return_tensors='pt', add_special_tokens=False).to('cpu')
    # print(sentence, inputs)  # Output: ['<s>'] {'input_ids': tensor([[0]], device='cpu')}
    _ = model(**inputs)      # After inputting the special token into the model, each layer of FFL will call register_forward_hook to store the activated neurons

outputs = torch.stack(outputs)    # output.shape=[12,1,102,3072]:layer=12, batch, promptLen+specialToken=102, neuronNum=3072
outputs = outputs[:, :, 1, :]    # 12 layers, 1, The dimension corresponding to the 1 token, neuron num = 3072
outputs = outputs.flatten()
```



# Activated Neurons

- Activated neurons in each layers
  - Input: ['realization', 'surprise', ..., 'remorse']

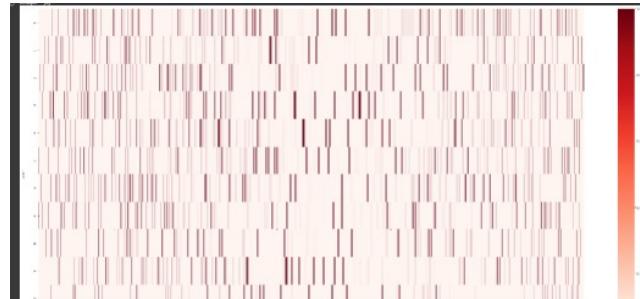
## 2.1 Input a prompt to show activated neuron

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

prompt = input("Please enter a prompt: ") # joy
print('Prompt: ', prompt)
df = pd.read_csv(f'./active_neuron_after_relu_csv/{prompt}-1.csv')

activatedNeuron = df.values
activatedNeuron = activatedNeuron.reshape(12,3072)
print(f'Activated Neurons{activatedNeuron.shape}: \n', activatedNeuron)

Please enter a prompt: joy
```





# Activated Neurons

- Cosine Similarity of Activated Neurons
  - Input: ['realization', 'surprise', ..., 'remorse']

## 2.2 Input two prompts to show cosine similarity of activated neurons

```
[ ] from sklearn.metrics.pairwise import cosine_similarity

prompt1 = input("Please enter prompt1: ") # joy
df1 = pd.read_csv(f'./active_neuron_after_relu_csv/{prompt1}-1.csv')
activatedNeuron1 = df1.values

prompt2 = input("Please enter prompt2: ") # love, admiration
df2 = pd.read_csv(f'./active_neuron_after_relu_csv/{prompt2}-1.csv')
activatedNeuron2 = df2.values

print('cosine similarity: ', cosine_similarity(activatedNeuron1, activatedNeuron2))

Please enter prompt1: joy
Please enter prompt2: love
cosine similarity:  [[0.49823117]]
```

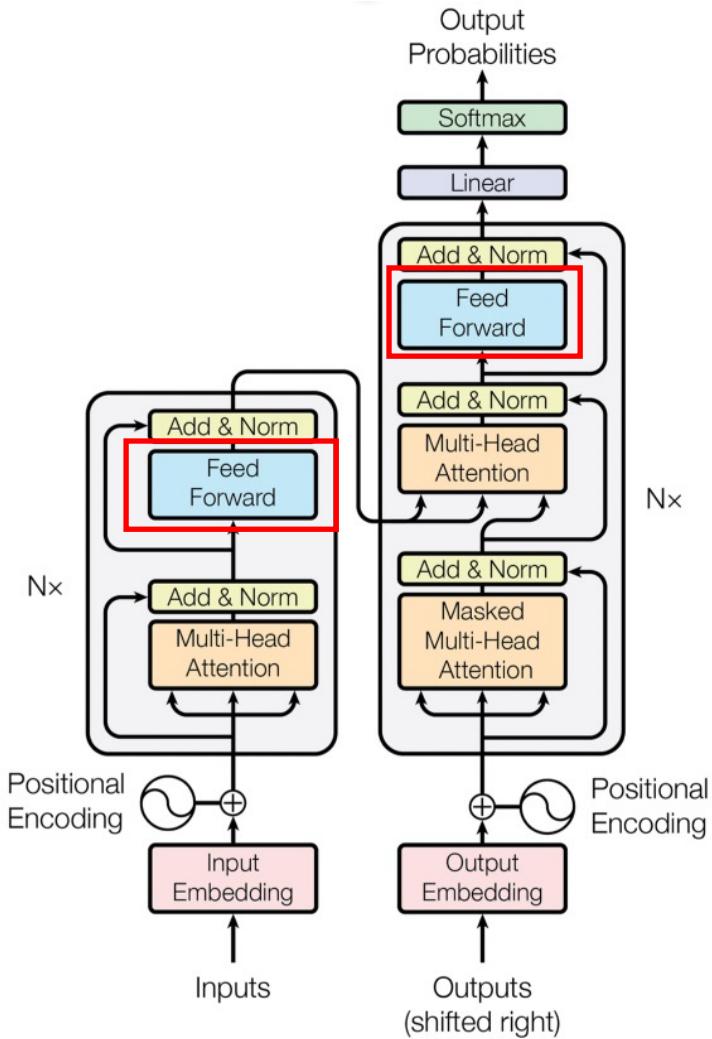


# Neurons in PLMs

THUNLP

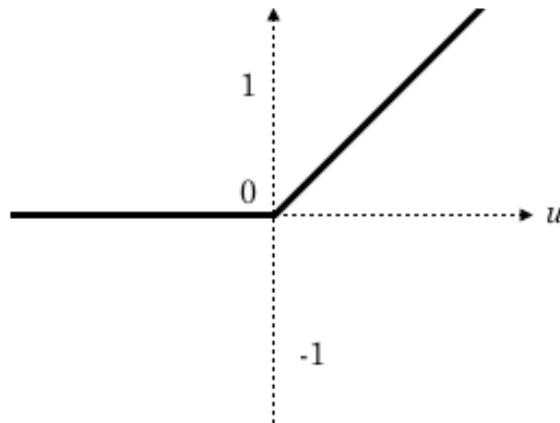


# Background: Neurons in FFNs

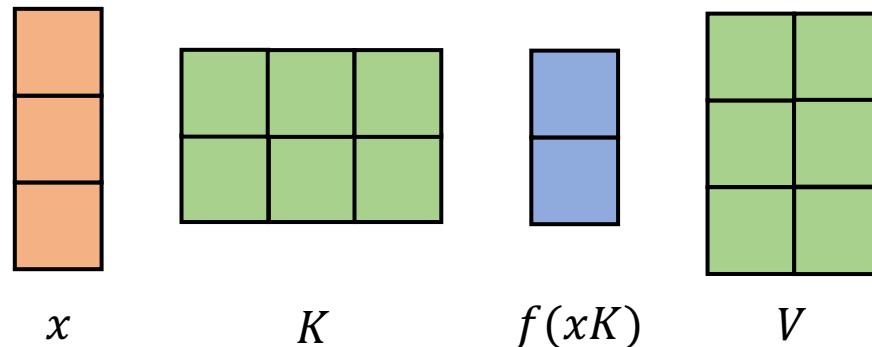


Transformer Architecture

$$f(u) = \max(0, u)$$



Activation Function

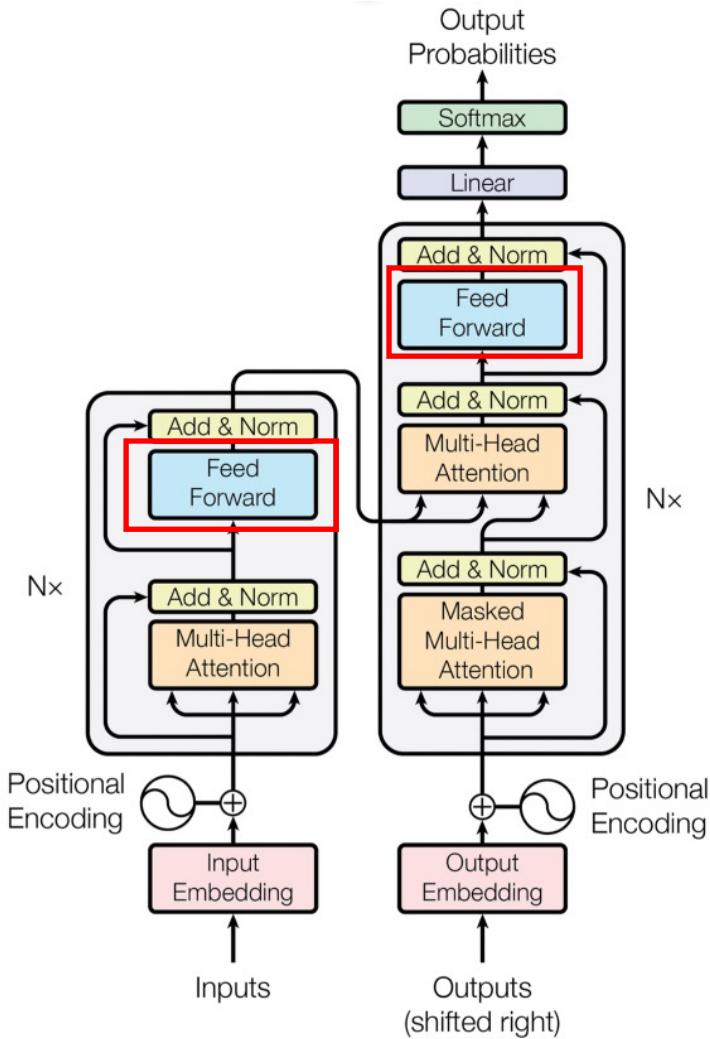


$$\text{FFN}(x) = f(\mathbf{x}\mathbf{K}^T + \mathbf{b}_1)\mathbf{V} + \mathbf{b}_2,$$

Feed Forward Neural Network

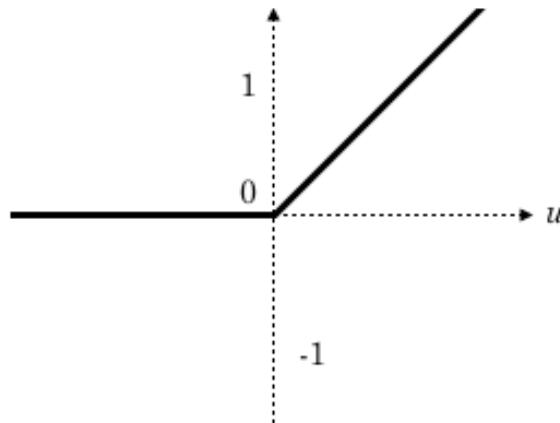


# Background: Neurons in FFNs

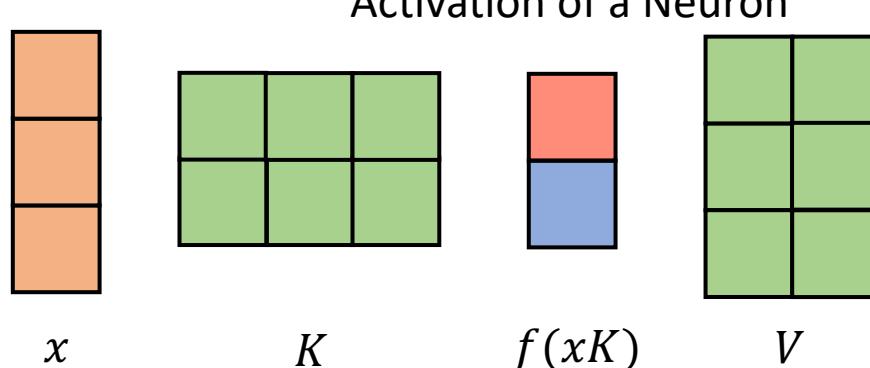


Transformer Architecture

$$f(u) = \max(0, u)$$



Activation Function

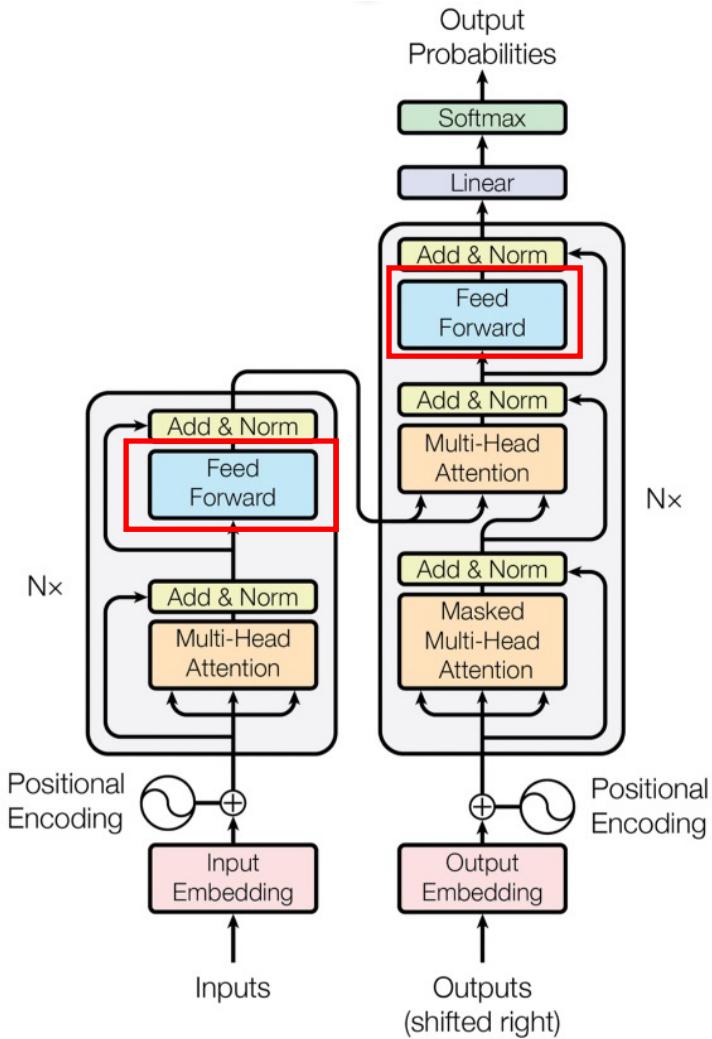


$$\text{FFN}(\mathbf{x}) = f(\mathbf{x}\mathbf{K}^\top + \mathbf{b}_1)\mathbf{V} + \mathbf{b}_2,$$

Feed Forward Neural Network

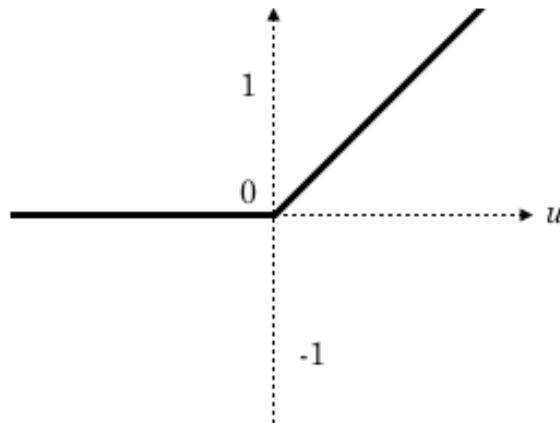


# Background: Neurons in FFNs



Transformer Architecture

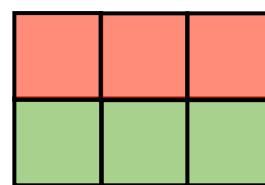
$$f(u) = \max(0, u)$$



Activation Function



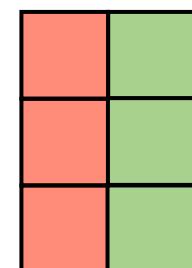
$x$



$K$



$f(xK)$



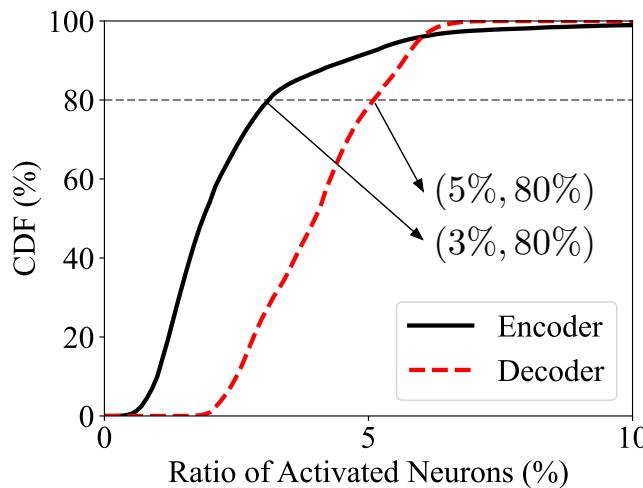
$V$

$$\text{FFN}(\mathbf{x}) = f(\mathbf{x}\mathbf{K}^\top + \mathbf{b}_1)\mathbf{V} + \mathbf{b}_2,$$

Feed Forward Neural Network

# Sparse activation phenomenon

- Sparse Activation Phenomenon in Large PLMs
- 80% inputs only activate less than 5% neurons of FFNs
- No useless neuron that keeps inactive for all inputs
- Related to Conditional Computation
  - Constrains a model to selectively activate parts of the neural network according to input



Cumulative distribution function (CDF) of the ratio of activated neurons in FFNs. Use T5-large (700 million parameters).



# Conditional computation

- Deep Learning of Representations: Looking Forward (Bengio, 2013)

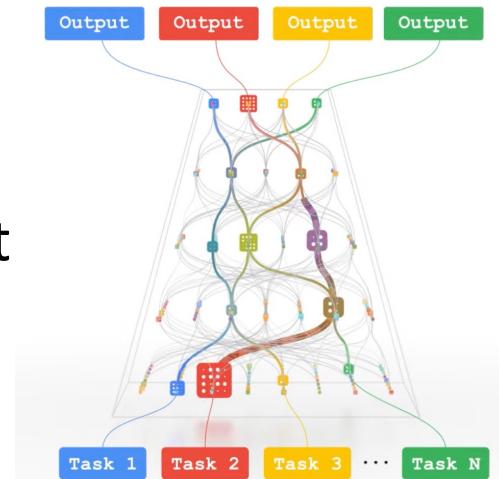
## 3 Scaling Computations

From a computation point of view, how do we scale the recent successes of deep learning to much larger models and huge datasets, such that the models are actually richer and capture a very large amount of information?

**Conditional Computation.** A central idea (that applies whether one parallelizes or not) that we put forward is that of *conditional computation*: instead of dropping out paths independently and at random, drop them in a learned and optimized way. Decision trees remain some of the most appealing

- Pathways (Jeff Dean, 2021)

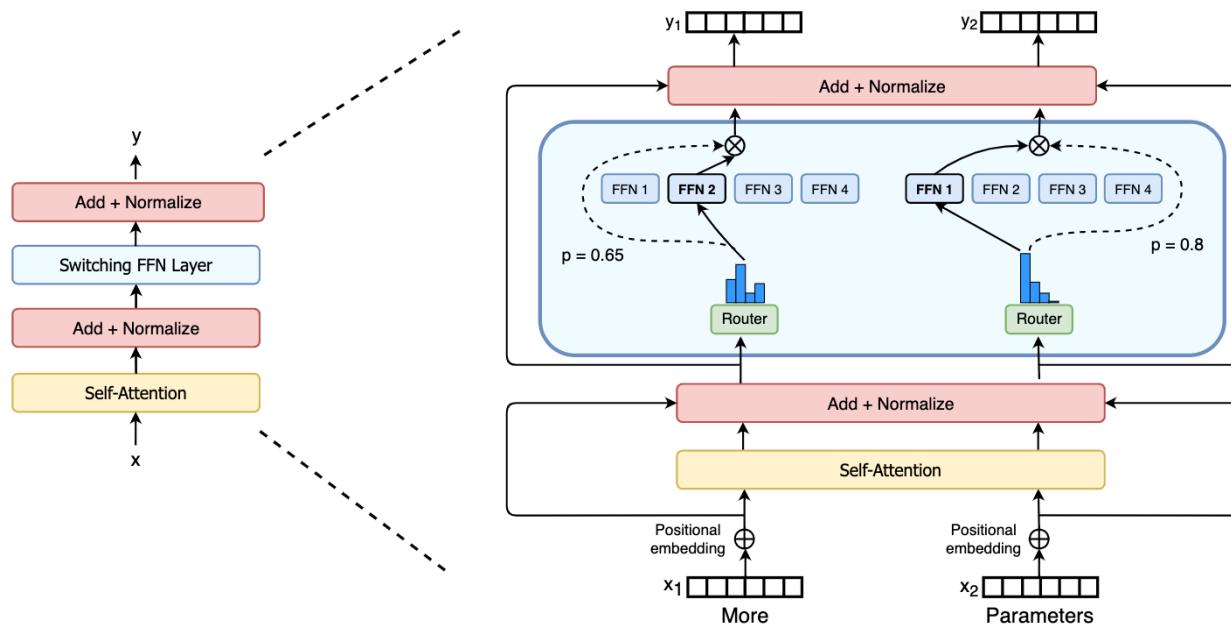
- Today's models are dense and inefficient
- Pathways will make them sparse and efficient





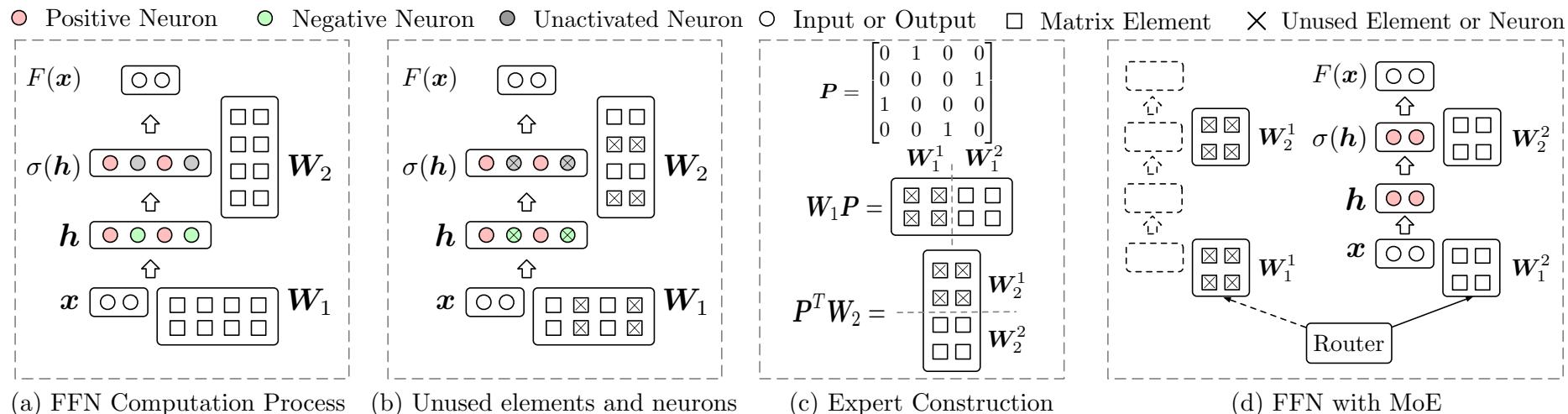
# MoEfication

- Mixture-of-experts (MoE)
- Use MoE to increase model parameters with tiny extra computational cost



# MoEification

- Split existing models into multiple experts while keeping model size unchanged





# MoEification

- Expert Construction
- Group the neurons that are often activated simultaneously
- Parameter Clustering Split
  - Treat the columns of  $W_1$  as a collection of vectors
  - K-means
- Co-Activation Graph Split
  - Construct a coactivation graph
  - Each neuron is represented by a node
  - Edge weight between two nodes is their co-activation value

$$\text{co-activation}(n, m) = \sum_{\mathbf{x}} \mathbf{h}_n^{(\mathbf{x})} \mathbf{h}_m^{(\mathbf{x})} \mathbb{1}_{\mathbf{h}_n^{(\mathbf{x})} > 0, \mathbf{h}_m^{(\mathbf{x})} > 0}$$



# MoEification

- Assign a score to each expert and select the experts with high scores

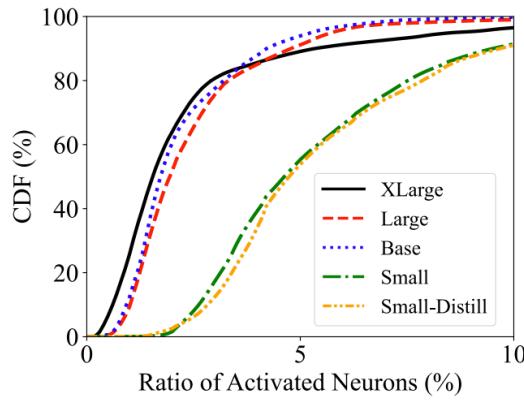
$$S = \arg \max_{A \subset \{1, 2, \dots, k\}, |A|=n} \sum_{i \in A} s_i$$

- Groundtruth Selection
  - Calculate the number of positive neurons in each expert as  $s_i$
- Parameter Center
  - Average all columns of  $W_1$  and use it as the center
- Learnable Router
  - Learn a router from the groundtruth on the training set

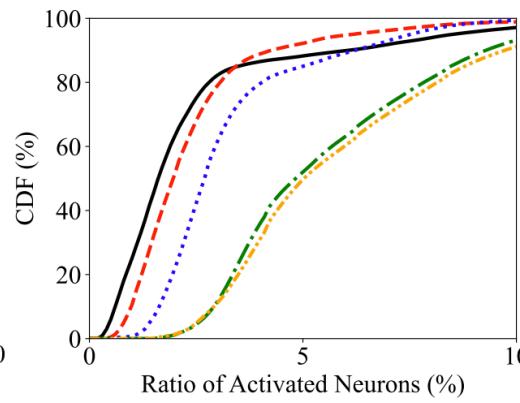


# MoEification

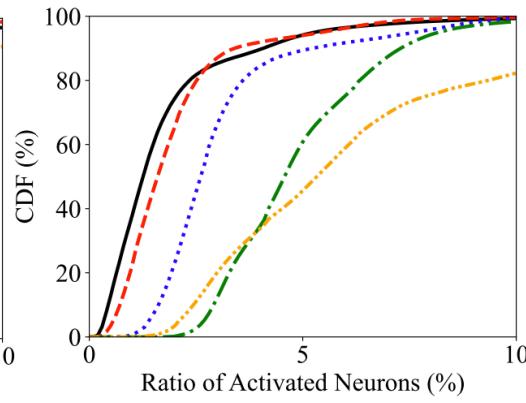
- Sparsity of Different T5 Models



(a) SST-2

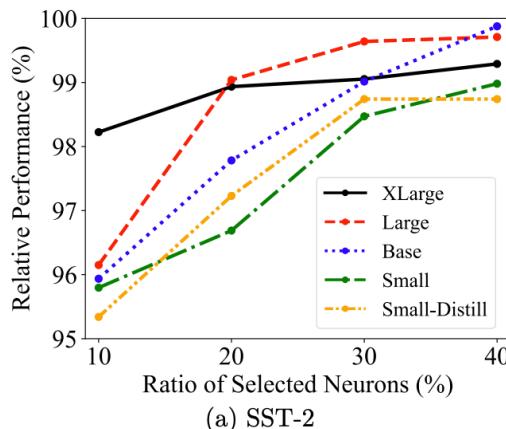


(b) MNLI

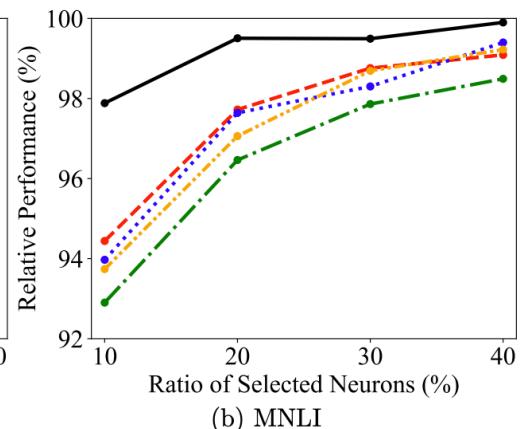


(c) RACE

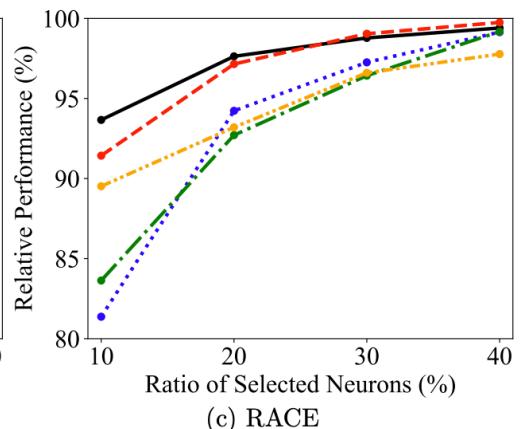
- MoEification with Different T5 Models



(a) SST-2



(b) MNLI



(c) RACE

# Observations on routing patterns

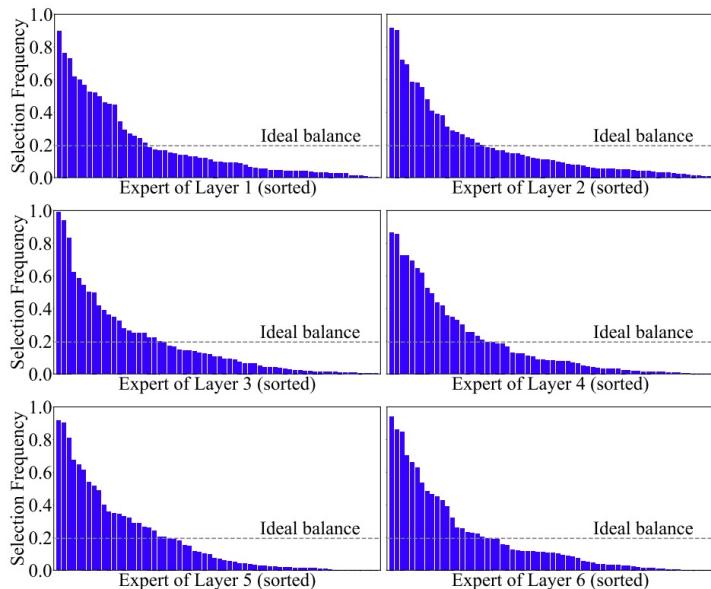
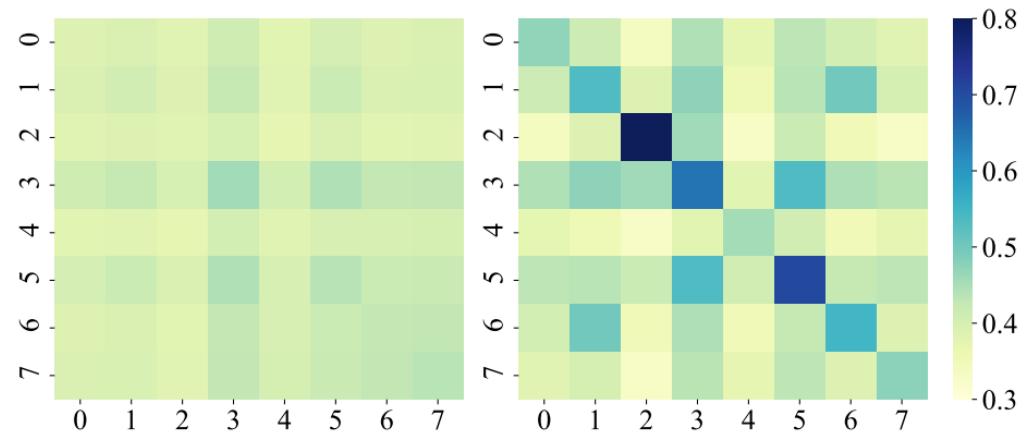


Figure 5: Selection Frequency of 64 experts in each encoder layer of MoEfied T5-Small. The frequency of ideal balance selection is 0.2 while the distribution is much unbalanced.



(a) The 8 most selected experts (b) The 8 least selected experts

Figure 6: Input similarities between experts in the last encoder layer of MoEfied T5-Small. For the most selected experts, both the self-similarities and inter-similarities are low. For the least selected experts, the self-similarities are much higher than inter-similarities.



# Analyze PLMs through neurons

- Specific Function
- Transferability Indicator
- Activated Neurons Can Reflect Human-Like Emotional Attributes

# Specific Function



# Expert units

- Identify whether the activation of a specific neuron can classify a concept
- $N_c^+$  positive sentences that contain concept c and  $N_c^-$  negative sentences that do not contain concept c.

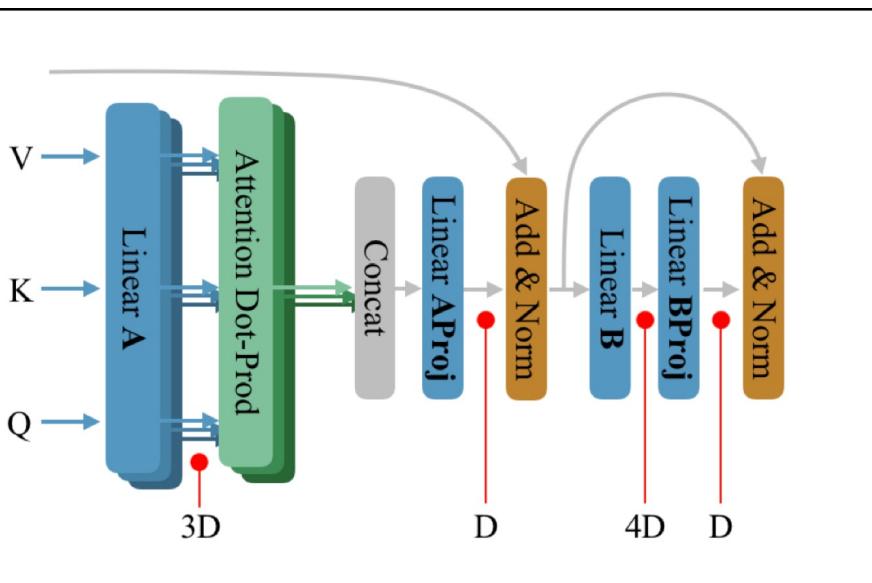
```
<instance docsrc="Indigenous architecture" id="shelter.00002">
  <answer instance="shelter.00002" senseid="shelter%1:06:00::" />
  <context>
    Types There are three traditional types of igloos ,
    all of different sizes and used for different purposes.
    The smallest were constructed as temporary
    <head>shelters</head>
    , usually only used for one or two nights .
  </context>
</instance>
```

## OneSec dataset

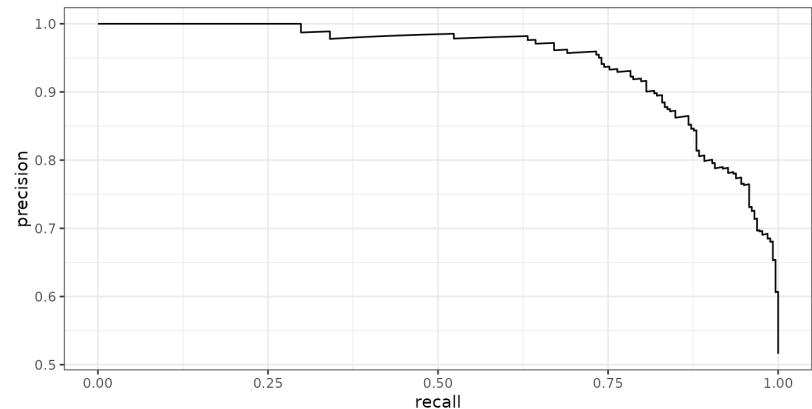
- **Sense:** Positive sentences contain a keyword with a specific WordNet sense, whereas negative sentences do not contain the keyword.
- **Homograph:** Positive sentences contain a keyword with a specific WordNet sense, whereas negative sentences contain the same keyword with a different meaning. Intuitively, *homograph* concepts are harder to disambiguate than *sense* concepts.

# Concept expertise

- Give each unit an index  $m$  and treat a unit as a binary classifier for the input sentences to compute AP



$$AP = \sum_n (R_n - R_{n-1}) P_n$$

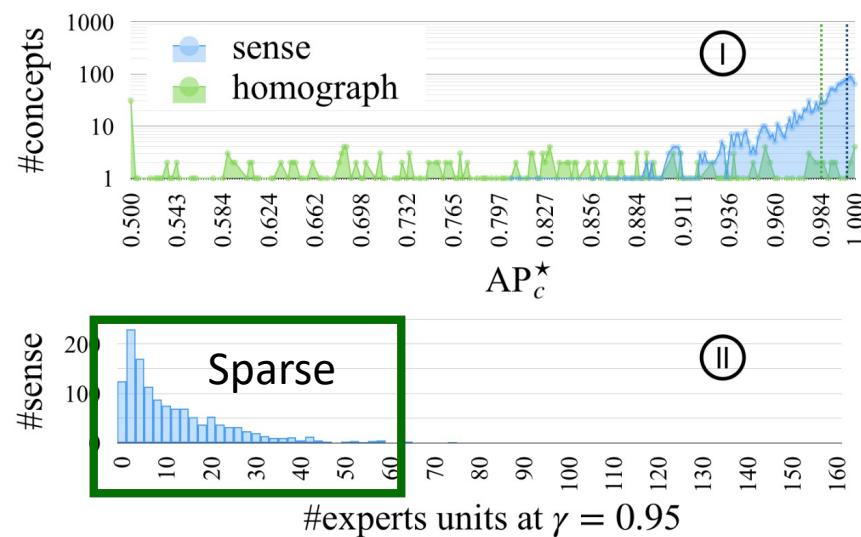
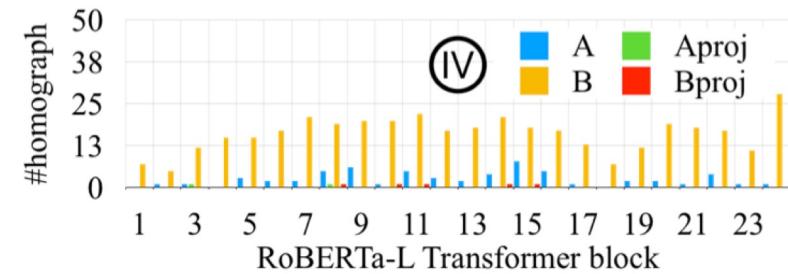
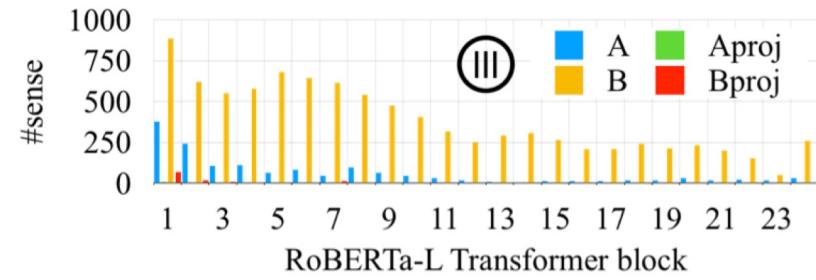
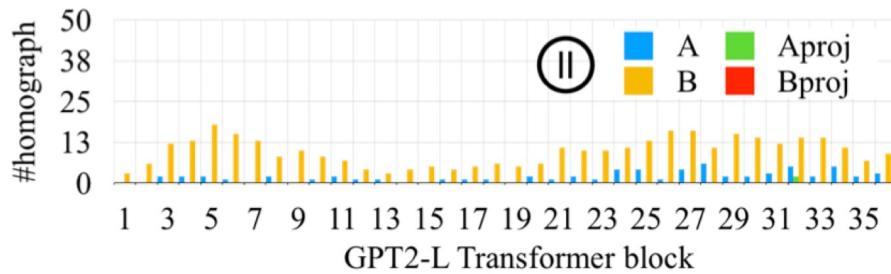
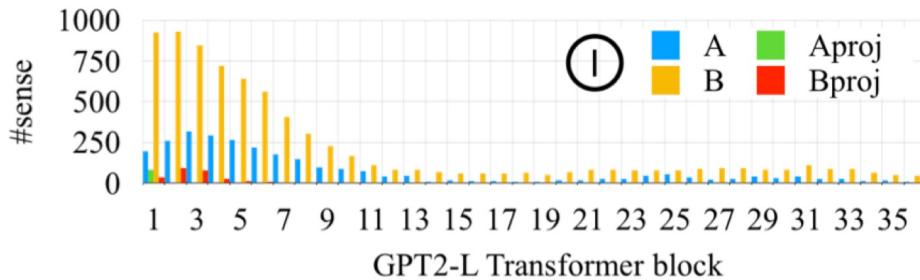


$$AP_c^* = \max_m \{ AP_c^m \}$$

Concept expertise

$$\chi_\gamma = \frac{|\{c \text{ s.t. } AP_c^* \geq \boxed{\gamma}\}|}{|C|} \quad \forall c \in C.$$

# Concept distribution





# Expertise and generalization results

- Detect the model's ability without fine-tuning

Models used	Task	(S)	(T)	$\mathcal{X}_{\gamma^*}$
BERT-B/L Distilbert RoBERTa-L XLM	GLUE Score	0.361	0.871	<b>0.892</b>
	CoLA	0.572	0.821	<b>0.874</b>
	SST-2	0.554	0.849	<b>0.864</b>
	MRPC (acc)	0.163	0.714	<b>0.753</b>
	MRPC (F1)	0.162	<b>0.916</b>	0.855
	STS-B (p)	0.000	<b>0.490</b>	0.401
	STS-B (s)	0.115	<b>0.903</b>	0.840
	QQP (acc)	0.340	0.936	<b>0.944</b>
	QQP (F1)	0.397	0.834	<b>0.859</b>
	MNLI-m	0.603	<b>0.833</b>	0.771
	MNLI-mm	0.452	0.906	<b>0.944</b>
	QNLI	0.286	0.857	<b>0.923</b>
	RTE	0.332	0.861	<b>0.873</b>
	WNLI	0.314	<b>0.751</b>	0.619
	AX	0.426	0.914	<b>0.956</b>
BERT-B/L DistilBERT RoBERTa-L	SQuAD 1.1 (F1)	<b>0.899</b>	0.840	0.850
	SQuAD 2.0 (F1)	<b>0.961</b>	0.752	0.937
	Average	0.408	0.826	<b>0.833</b>



# Concept overlap

- Let the overlap between concepts  $q$  and  $v$  be

$$\Omega(q, v) = \frac{\|\mathbf{s}_q \cap \mathbf{s}_v\|_1}{\|\mathbf{s}_q \cup \mathbf{s}_v\|_1} \in [0, 1]$$

Query definition	Concept	$\Omega(q, v)$
A seat for one person, with a support for the back.	chair%1:06:00 (query)	1.000
	table%1:06:01	0.458
	bed%1:06:00	0.361
	cup%1:06:00	0.341
	table%1:06:01 VS. table%1:14:00	0.336
	floor%1:06:00	0.328
The position of professor.	chair%1:04:00 (query)	1.000
	chair%1:04:00 VS. chair%1:06:00	0.575
	fellow%1:18:02	0.371
	director%1:18:03	0.297
	administration%1:04:00	0.243
	member%1:18:00	0.241



# Conditioned text generation

- Selected expert units to compute

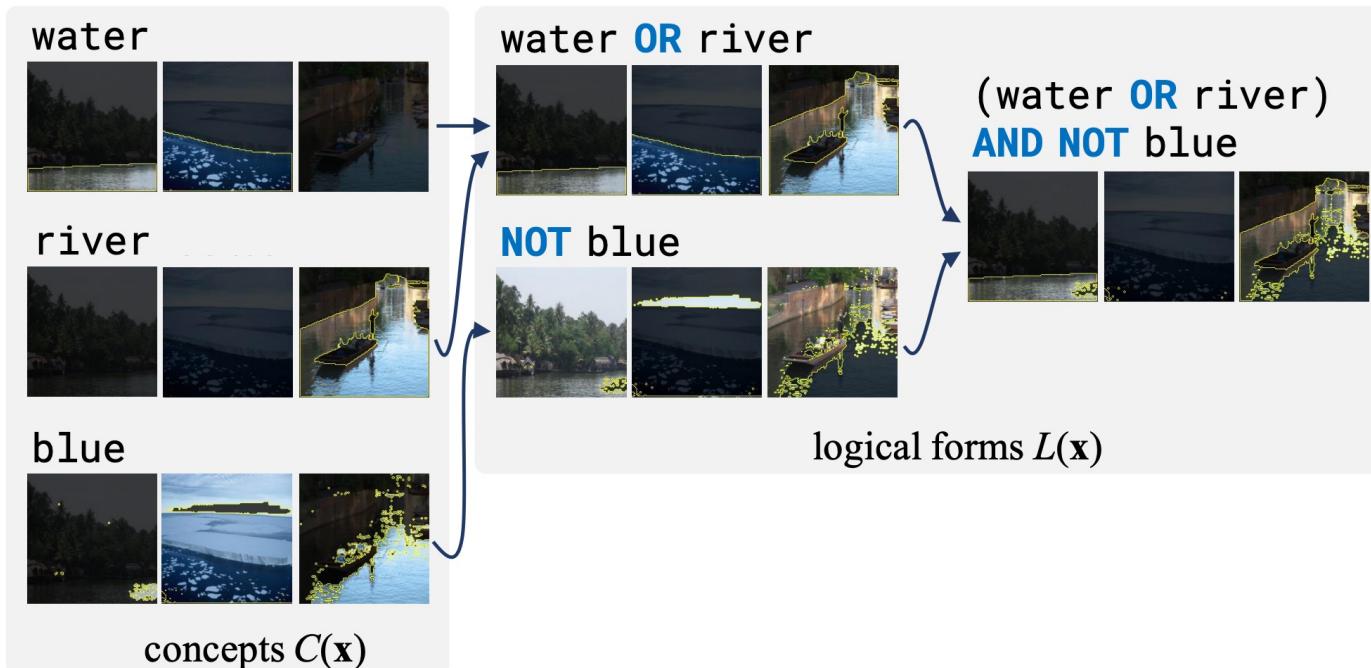
K forced	Once upon a time + Generated induced for concept bird%1:05:00 (warm-blooded egg-laying vertebrates)
0 (0%)	, I had a friend who used to teach high school English and he was like, "Oh, all you have to do is just get out there
40 (0.009%)	, many of these treasures were worth hundreds of thousands of dollars.\n But this isn't the first time that a horse has been
60 (0.015%)	, through a freak occurrence, an invasion of house sparrows which so often reduces the black-browed this nation recreates through
80 (0.019%)	, our own ancestors rode about on chicken-like air wings.\n But this wonder of the air has no such wings.\n Taking down
200 (0.048%)	of year, birds chase each and watching. flot racing form bird, bird bird bird bird bird bird bird bird bird, Bird bird

Once upon a time + Generated induced for concept lead%1:07:02 (an advantage held by a competitor in a race)
50 (0.012%) the left-hander would always start at the front in the first two instances, but when Mauricio Gaponi rose to the podium,

Once upon a time + Generated induced for concept lead%1:27:00 (a soft heavy toxic malleable metallic element)
100 (0.024%) a crust layer was applied to a partially fortified nickel base, thereby causing to zinc- and copper- ground element cob. The occurrence of those metal and chrome

# Compositional explanations of neurons

- Neurons learn compositional concepts



- Compositional explanations allow users to predictably manipulate model behavior

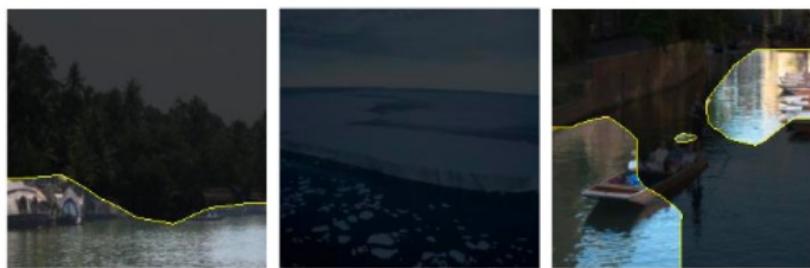


# Find neurons

- For an individual neuron, thresholding its activation



(b) neuron  $f_{483}(\mathbf{x})$

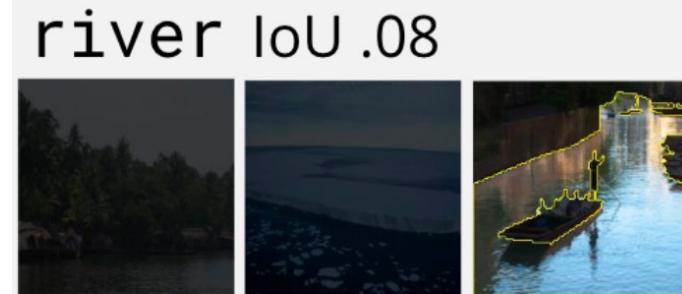
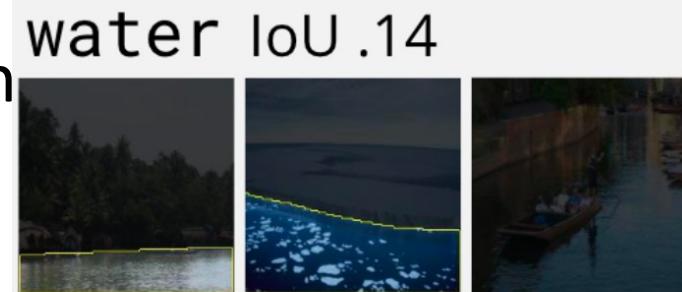


(c) neuron masks  $M_{483}(\mathbf{x})$



# Find neurons

- For an individual neuron, thresholding
- Compare with the mask of concepts



(d) concepts  $C(\mathbf{x})$



# Find neurons

- Search for the most similar concept

$$\text{EXPLAIN-NETDISSECT}(n) = \arg \max_{C \in \mathcal{C}} \delta(n, C)$$

$$\delta(n, C) \triangleq \text{IoU}(n, C) = \left[ \sum_{\mathbf{x}} \mathbb{1}(M_n(\mathbf{x}) \wedge C(\mathbf{x})) \right] / \left[ \sum_{\mathbf{x}} \mathbb{1}(M_n(\mathbf{x}) \vee C(\mathbf{x})) \right]$$

Find logical forms induced from the concepts

Compose these concepts via compositional operations: AND OR NOT



# Tasks

- Image Classification
  - Scene Recognition
  - ResNet-18
- NLI
  - SNLI
    - BiLSTM+MLP
    - Probe neurons in MLP, input is premise-hypothesis pairs
  - Concepts:
    - Penn Treebank POS tags + 2000 most common words
      - Appear in premise or hypothesis
    - Whether premise and hypothesis have more than 0%, 25%, 50%, or 75% word overlap
  - Additional Operator
    - NEIGHBORS( $C$ ), the union of 5 most close words with  $C$ 
      - Judged by cosine similarity of Glove embeddings



# Do neurons learn compositional concepts?

## Unit 870 (gender-sensitive)

((((NOT hyp:man) AND pre:man) OR hyp:eating)  
AND (NOT pre:woman)) OR hyp:dancing  
IoU **0.123** W<sub>entail</sub> **-0.046** W<sub>neutral</sub> **-0.021** W<sub>contra</sub> **0.040**

**Pre** A guy pointing at a giant blackberry.  
**Hyp** A woman tearing down a giant display.  
Act **29.31** True **contra** Pred **contra**  
**Pre** A man in a hat is working with...flowers.  
**Hyp** Women are working with flowers.  
Act **27.64** True **contra** Pred **contra**

## Unit 99 (high overlap)

((NOT hyp:JJ) AND overlap-75% AND (NOT  
pre:people)) OR pre:basket OR pre:tv  
IoU **0.118** W<sub>entail</sub> **0.043** W<sub>neutral</sub> **-0.029** W<sub>contra</sub> **-0.021**

**Pre** A woman in a light blue jacket is riding a bike.  
**Hyp** A women in a jacket riding a bike.  
Act **19.13** True **entail** Pred **entail**  
**Pre** A girl in a pumpkin dress sitting at a table.  
**Hyp** There is a girl in a pumpkin dress sitting at a table.  
Act **17.84** True **entail** Pred **entail**

## Unit 15 (sitting only in hypothesis)

hyp:eating OR hyp:sitting OR hyp:sleeping  
OR hyp:sits AND (NOT pre:sits)  
IoU **0.239** W<sub>entail</sub> **-0.083** W<sub>neutral</sub> **-0.059** W<sub>contra</sub> **0.086**

**Pre** A person...is walking through an airport.  
**Hyp** A woman sits in the lobby waiting on the doctor.  
Act **30.68** True **contra** Pred **contra**  
**Pre** A man jumps over another man...  
**Hyp** Two men are sitting down, watching the game.  
Act **27.64** True **contra** Pred **contra**

## Unit 473 (unclear)

((NOT hyp:sleeping) AND (pre>NN OR pre:NNS))  
AND (NOT hyp:alone) AND (NOT hyp:nobody)  
IoU **0.586** W<sub>entail</sub> **0.020** W<sub>neutral</sub> **0.016** W<sub>contra</sub> **-0.050**

**Pre** A gentleman in a striped shirt gesturing with a stick...  
**Hyp** A gentleman in a striped shirt joyously gesturing.  
Act **31.62** True **neutral** Pred **neutral**  
**Pre** An Asian man in a...uniform diving...in a game.  
**Hyp** A person in a uniform does something.  
Act **29.76** True **neutral** Pred **entail**



# Can we target explanations to change model behavior?

## Unit 39 (nobody in hypothesis)

hyp:nobody AND (NOT pre:hair) AND (NOT pre:RB) AND (NOT pre:’s)  
IoU **0.465** W<sub>entail</sub> -0.117 W<sub>neutral</sub> -0.053 W<sub>contra</sub> 0.047

**Pre** Three women prepare a meal in a kitchen.  
**Orig Hyp** The ladies are cooking.  
**Adv Hyp** **Nobody but** the ladies are cooking.  
True entail <sup>adv</sup>→ neutral Pred entail <sup>adv</sup>→ contra

## Unit 15 (sitting only in hypothesis)

hyp:eating OR hyp:sitting OR hyp:sleeping OR  
hyp:sits AND (NOT pre:sits)  
IoU **0.239** W<sub>entail</sub> -0.083 W<sub>neutral</sub> -0.059 W<sub>contra</sub> 0.086

**Orig Pre** A blond woman is holding 2 golf balls while reaching down into a golf hole.  
**Adv Pre** A blond woman is holding 2 golf balls.  
**Hyp** A blond woman is sitting down.  
True contra <sup>adv</sup>→ neutral Pred contra <sup>adv</sup>→ contra

## Unit 133 (couch words in hypothesis)

NEIGHBORS(hyp:couch) OR hyp:inside OR  
hyp:home OR hyp:indoors OR hype:eating  
IoU **0.202** W<sub>entail</sub> -0.125 W<sub>neutral</sub> -0.024 W<sub>contra</sub> 0.088

**Pre** 5 women sit around a table doing some crafts.  
**Orig Hyp** 5 women sit around a table.  
**Adv Hyp** 5 women sit around a table **near a couch**.  
True entail <sup>adv</sup>→ neutral Pred entail <sup>adv</sup>→ contra

## Unit 941 (inside/indoors in hypothesis)

hyp:inside OR hyp:not OR hyp:indoors OR  
hyp:moving OR hyp:something  
IoU **0.151** W<sub>entail</sub> 0.086 W<sub>neutral</sub> -0.030 W<sub>contra</sub> -0.023

**Orig Pre** Two people are sitting in a station.  
**Adv Pre** Two people are sitting in a **pool**.  
**Hyp** A couple of people are inside and not standing.  
True entail <sup>adv</sup>→ neutral Pred entail <sup>adv</sup>→ entail

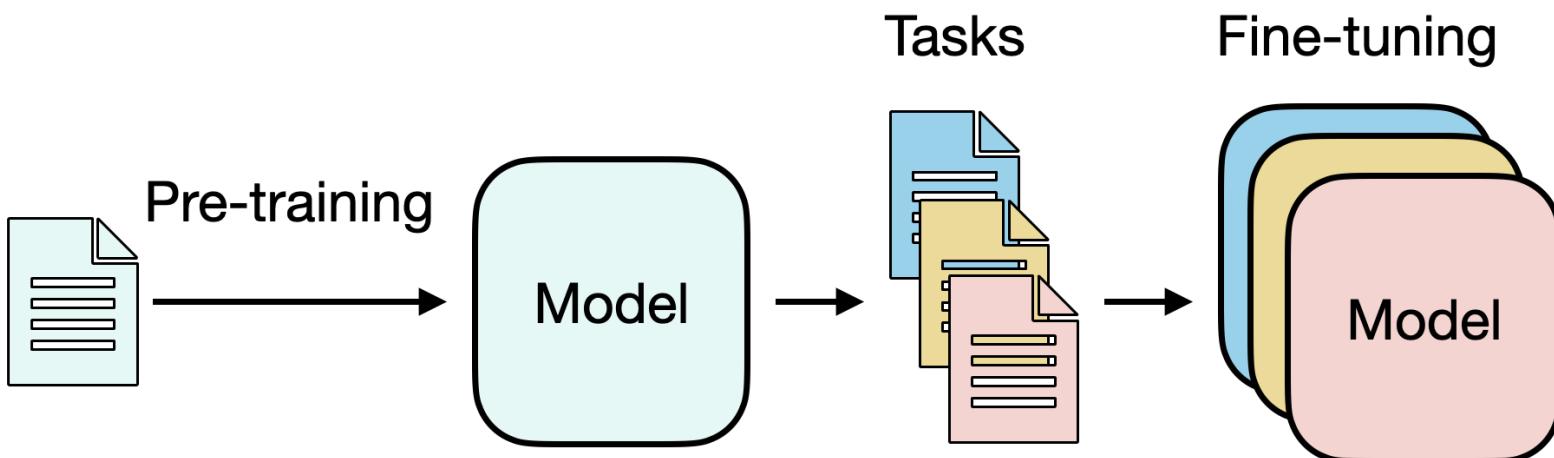


# Cognitive Abilities of Big Models

THUNLP

# Task generalization of PLM

- Question: why can PLMs easily adapt to various NLP tasks even with **small-scale** data?





# Task generalization of PLM

- PLM acquires versatile **knowledge** during pre-training, which can be leveraged to solve various tasks



Customer Service  
ActiveChat.ai



Games  
AI Dungeon



Chatbots  
AI Buddy



Customer Service  
Chatdesk



LegalTech  
aiLawDocs



Semantic Search  
Algolia



Chatbots  
AskBrian

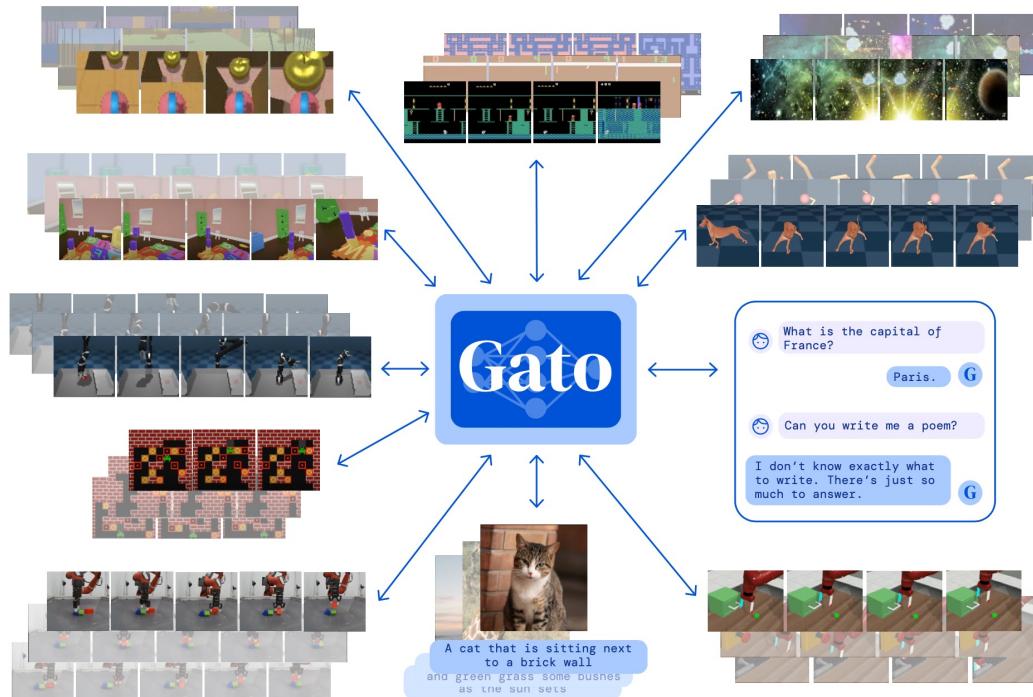


AI Copywriting  
Conto AI



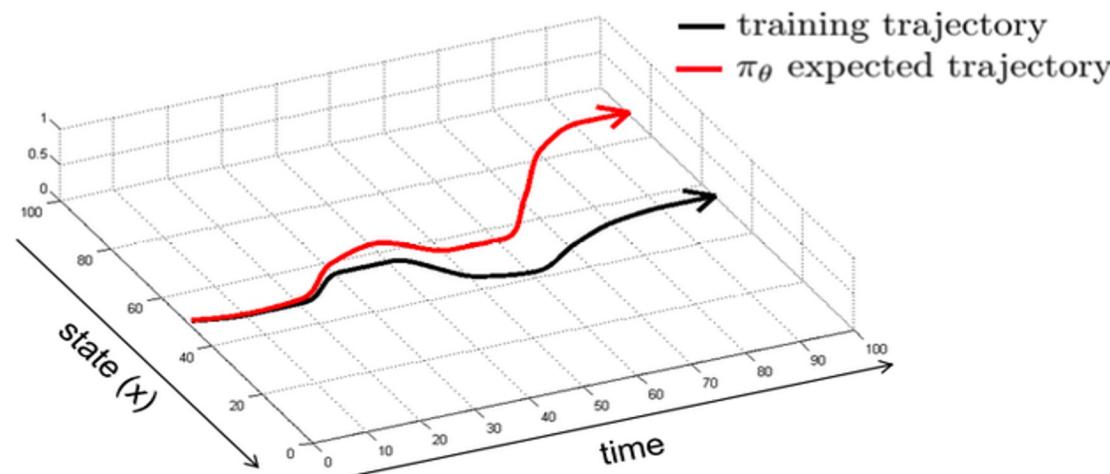
# Cognitive abilities of PLMs

- Recent studies have shown that PLMs also have **cognitive abilities**, and can manipulate existing tools to complete a series of complex tasks



# Fundamentals & framework

- Imitation learning in RL
  - Learning from behaviors instead of accumulating rewards
  - State-action pairs  $\mathcal{D} = \{(s_1, a_1), (s_2, a_2), (s_3, a_3), \dots\}$
  - State as features and action as labels
  - Target: imitate the trajectory of behaviors





# Fundamentals & framework

- Large-scale pre-trained models
  - Universal knowledge learned from pre-training
- Interactive space
  - An **environment** that models could interact with
  - **State space:** display states and changes
  - **Action space:** a set of pre-defined actions in the environment



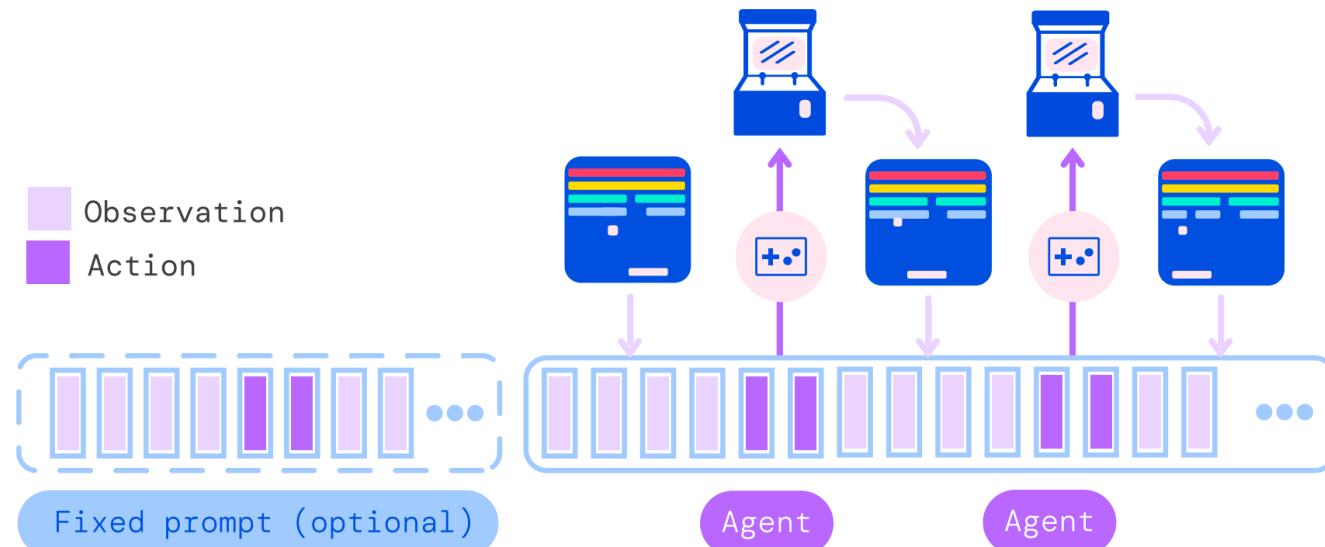
# Fundamentals & framework

- Given a goal, we model each action and state to achieve the goal in a unified space by a PLM
- Tokenization
  - Tokenize human behaviors (**actions** in the action space) and **states** in the state space to a same space
  - The tokenized information could be processed by PLM
- Directly training
  - The behaviors could be **autoregressively** predicted



# Outline

- Interactive space:
  - Search engine
  - WebShop
  - Sandbox





# Interactive space: search engine

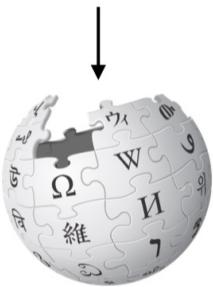
- A rising challenge in NLP is long-form QA
  - A paragraph-length answer is generated in response to an open-ended question

<b>Question</b>	Why did we decide that certain words were "bad" and shouldn't be used in social settings?
<b>Answer</b>	We have all been taught that there are certain words that should not be used in public, at least not without good reason. But why do some words fall into this category while others do not? It turns out that there is not a single answer to this question, as each word has a different history and carries different meanings in different cultures and families [2]. However, there are some general categories that "bad" words tend to fall into. "Bad" words generally relate to parts of life that we don't like talking about in public, like bathroom functions, or negative ways of talking about people's religion, skin color, ability, or gender [3]. Often, people say bad words when they are experiencing strong emotions, and those words are linked to our emotions [3].
<b>References (titles only)</b>	[1, 2, 3] Why Are Some Words 'Bad'?   Vermont Public Radio ( <a href="http://www.vpr.org">www.vpr.org</a> ) [4] On Words: 'Bad' Words and Why We Should Study Them   UVA Today ( <a href="http://news.virginia.edu">news.virginia.edu</a> ) [5] The Science of Curse Words: Why The &@#! Do We Swear? ( <a href="http://www.babbel.com">www.babbel.com</a> )

# Interactive space: search engine

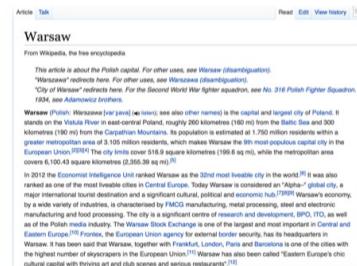
- The task has two core components:
  - information retrieval and information synthesis

Q: How many of Warsaw's inhabitants spoke Polish in 1933?



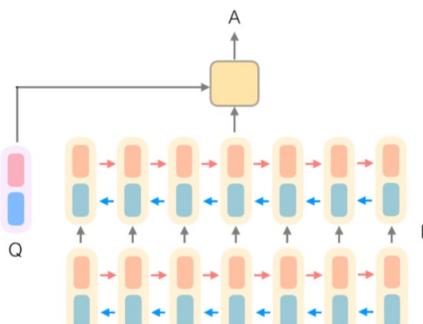
**WIKIPEDIA**  
The Free Encyclopedia

**Document  
Retriever**



**Document  
Reader**

833,500





# Interactive space: search engine

- WebGPT
  - Outsource document retrieval to the Microsoft Bing Web Search API
  - Utilize unsupervised pre-training to achieve high-quality document synthesis by fine-tuning GPT-3
  - Create a text-based web-browsing environment that both humans and language models can interact with



# Interactive space: search engine

- Text-based web-browser

How can I train the crows in my neighborhood to bring me gifts?

This question does not make sense   This question should not be answered

Search results for: how to train crows to bring you gifts   Quotes ↗

← how to train crows to bring Find in page + Add new quote

[How to Make Friends With Crows - PetHelpful](#)  
If you did this a few times, your crows would learn your new place, but as I said, I'm not sure if they will follow or visit you there since it's probably not in their territory. The other option is simply to make new crow friends with the crows that live in your new neighborhood.

[Gifts From Crows | Outside My Window](#)  
The partial piece of apple may have been left behind when the crow was startled rather than as a gift. If the crows bring bright objects you'll know for sure that it's a gift because it's not something they eat. Brandi Williams says: May 28, 2020 at 7:19 am.

[1] Gifts From Crows | Outside My Window  
([www.birdsoutsidemywindow.org](http://www.birdsoutsidemywindow.org))

Many animals give gifts to members of their own species but crows and other corvids are the only ones known to give gifts to humans.

Number of quote tokens left: 463  
Number of actions left: 96

Done quoting! Write an answer



# Interactive space: search engine

- Text-based web-browser

Command	Effect
Search <query>	Send <query> to the Bing API and display a search results page
Clicked on link <link ID>	Follow the link with the given ID to a new page
Find in page: <text>	Find the next occurrence of <text> and scroll to it
Quote: <text>	If <text> is found in the current page, add it as a reference
Scrolled down <1, 2, 3>	Scroll down a number of times
Scrolled up <1, 2, 3>	Scroll up a number of times
Top	Scroll to the top of the page
Back	Go to the previous page
End: Answer	End browsing and move to answering phase
End: <Nonsense, Controversial>	End browsing and skip answering phase



# Interactive space: search engine

- Text-based web-browser

- ◆ Question

- How can I train the crows in my neighborhood to bring me gifts?

- ◆ Quotes

- From Gifts From Crows | Outside My Window ([www.birdsoutsidemywindow.org](http://www.birdsoutsidemywindow.org))

- > Many animals give gifts to members of their own species but crows and other corvids are the only ones known to give gifts to humans.

- ◆ Past actions

- Search how to train crows to bring you gifts

- Click Gifts From Crows | Outside My Window [www.birdsoutsidemywindow.org](http://www.birdsoutsidemywindow.org)

- Quote

- Back

- ◆ Title

- Search results for: how to train crows to bring you gifts

- ◆ Scrollbar: 0 - 11

- ◆ Text

- [0] How to Make Friends With Crows - PetHelpful [pethelpful.com]

- If you did this a few times, your crows would learn your new place, but as I said, I'm not sure if they will follow or visit you there since it's probably not in their territory. The other option is simply to make new crow friends with the crows that live in your new neighborhood.

- [1] Gifts From Crows | Outside My Window [www.birdsoutsidemywindow.org]

- The partial piece of apple may have been left behind when the crow was startled rather than as a gift. If the crows bring bright objects you'll know for sure that it's a gift because it's not something they eat.

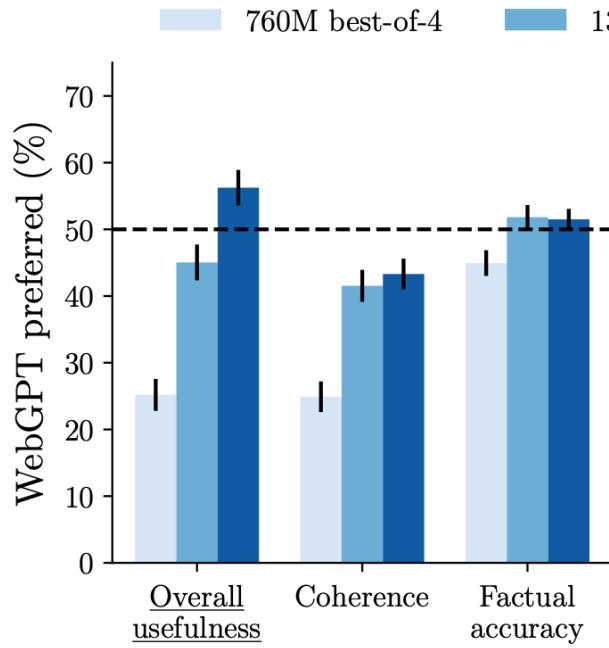
- Brandi Williams says: May 28, 2020 at 7:19 am.

- ◆ Actions left: 96

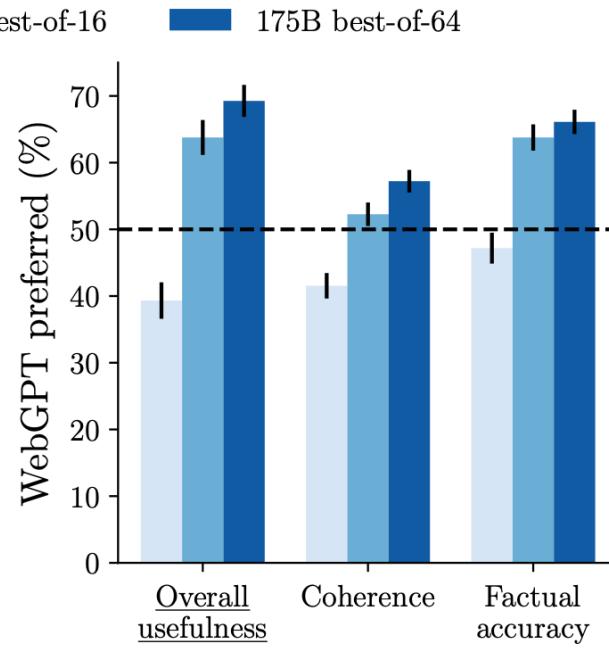
- ◆ Next action

# Interactive space: search engine

- WebGPT-produced answers are **more preferred** than human-generated ones
- Better **coherence** and **factual accuracy**



(a) WebGPT vs. human demonstrations.



(b) WebGPT vs. ELI5 reference answers.



# Interactive space: search engine

- An example:
  - How does neural networks work?

how do neural networks work

```
.d88888b.  
.8P"      "9bd888b.  
.8P      .d8P"  '"988.  
.8888     .d8P" ,   98.  
.8P" 88    8"      .d98b.  88  
.8P    88    8 .d8P"  "98b. 88  
88    88    8P"  '"8b.    "98.  
88.   88    8      8"8b.  88  
88    "98.8   8    88   "88  
'8b.   "98., .d8   88   88  
88   "98b.   .d8P" 8    88   d8"  
88    "98bP"   .8    88   .d8"  
"8b    '      .d8P"  8888"  
"88b.,   .d8P"   d8"  
"9888P98b.   .d8"  
"988888P"
```



# Interactive space: WebShop

- WebShop for online shopping (HTML mode)

**A**

**WebShop search**

**Instruction:**  
I'm looking for a small portable folding desk that is already fully assembled; it should have a khaki wood finish, and price lower than 140.00 dollars

portable folding desk khaki wood **1** **Search**

**Back to Search** **2 results**

Page 1 (Total results: 50) **Next >**

**B09Q3B186B**  
MENHG Folding Breakfast Tray Table, Efficient Home Laptop Notebook Computer Desk, Portable Writing Study Desk, Sturdy Home Office Table Workstation  
\$109.0

**B09P5ZBCWR**  
KPSP Folding Study Desk Bed Breakfast Serving Tray Table Efficient Home Laptop Notebook Computer Desk Portable Standing Desk for Small Space Bedroom

Description: Product laptop desk. Product material: walnut. Product weight: 4.6 pounds. Material: high quality thick steel pipe, black brushed sheet. Special design: black brushed smooth table top, increase the length and width of the table, it is possible to place the computer and various items. Function: Can be used as computer desk, dining table, bedside table. Product size: 23.6x15.7x11 inches

• **item-detail** 【Large Size】 Styling with light wood. Holds laptops up to 17 inches. It also have spacious space (23.6x15.7x11 inches) for your laptop, notebook, mouse, pen and coffee. Its generous size gives this versatile desk even more flexibility.

• **item-detail** 【Wide Application】 Our foldable lap desk can be used as a



MENHG Folding Laptop Table Bed Desk PC Lap Desk with Drawer Book Stand Reading Holder Leg Space Laptop Bed Tray Foldable Lazy Table Breakfast Desk Sofa Small Desk for Small Space  
Price: \$109.0  
Rating: N.A.

**Description** **Overview** **Buy Now** **5**

**Color** **black** **khaki** **white** **3**

**Reward:** **1.0**



# Interactive space: WebShop

- Simple mode which strips away extraneous meta-data from raw HTML into a simpler format

**B**  **Simple mode**

I'm looking for a small portable folding desk that is already fully assembled [...]

[btn] Back to Search [/btn]  
Page 1 (Total results: 50) [btn] Next [/btn]  
[btn] MENHG Folding Breakfast Tray [...] [/btn]  
\$109.0  
[btn] KPSP Folding Study Desk Bed [...] [/btn]

**C**

$\bar{u}$  (Instruction): I'm looking for a small portable...

$\bar{y}$  (Description): MENHG Folding Laptop Table Bed...

$y_{price}$ : \$109.0

$Y_{opt}$  (Options): { black, khaki, white }

$Y_{att}$  (Attributes): { steel pipe, no assembly, portable }



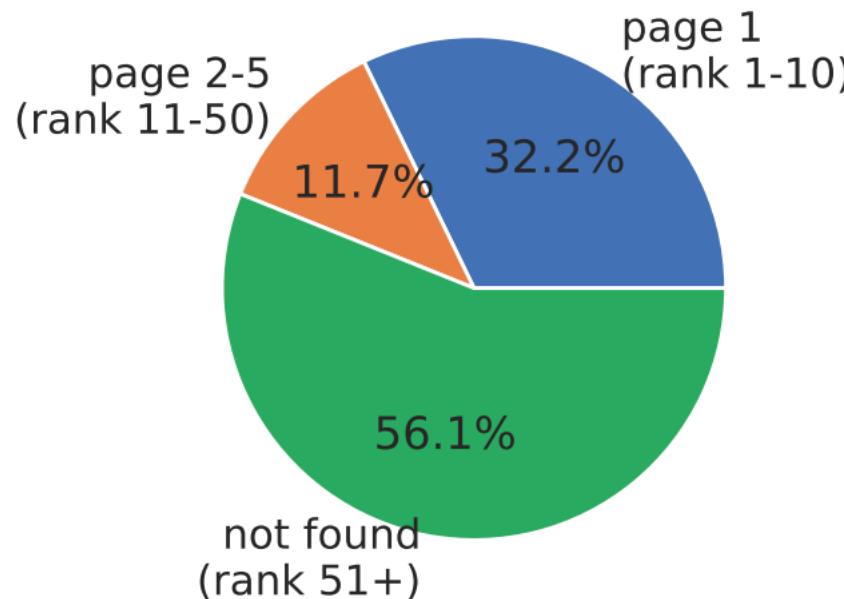
# Interactive space: WebShop

- Actions in WebShop

Type	Argument	State → Next State
search	[ <i>Query</i> ]	Search → Results
choose	Back to search	* → Search
choose	Prev/Next page	Results → Results
choose	[ <i>Product title</i> ]	Results → Item
choose	[ <i>Option</i> ]	Item → Item
choose	Desc/Overview	Item → Item-Detail
choose	Previous	Item-Detail → Item
choose	Buy	Item → Episode End

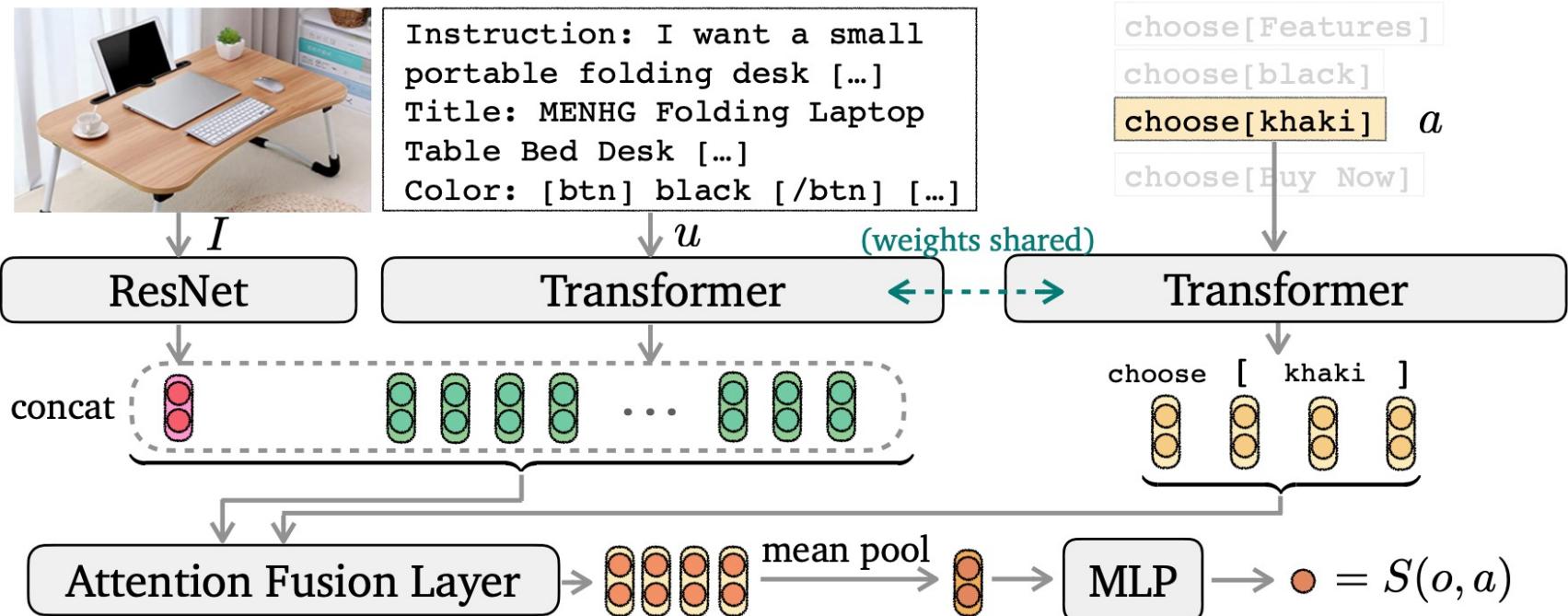
# Interactive space: WebShop

- Item rank in search results when the instruction is directly used as search query



# Interactive space: WebShop

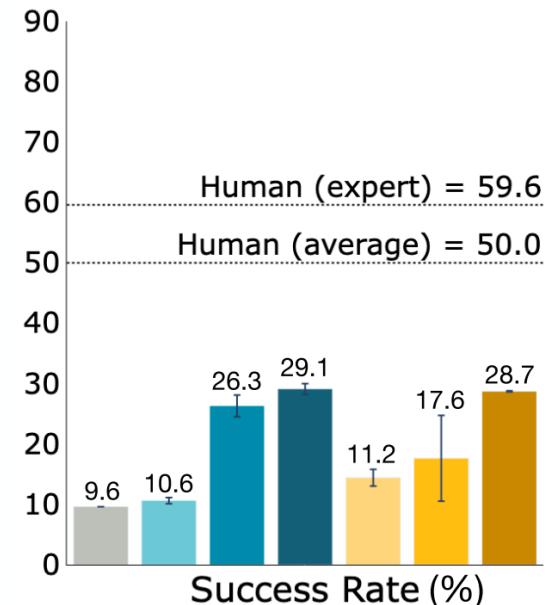
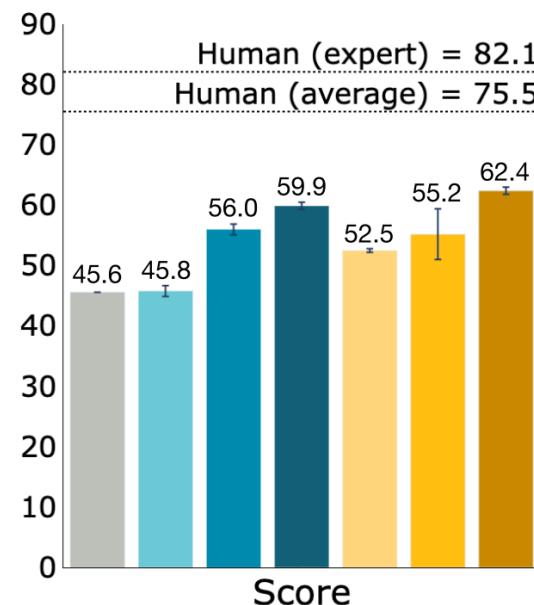
- Model implementation



# Interactive space: WebShop

- Results

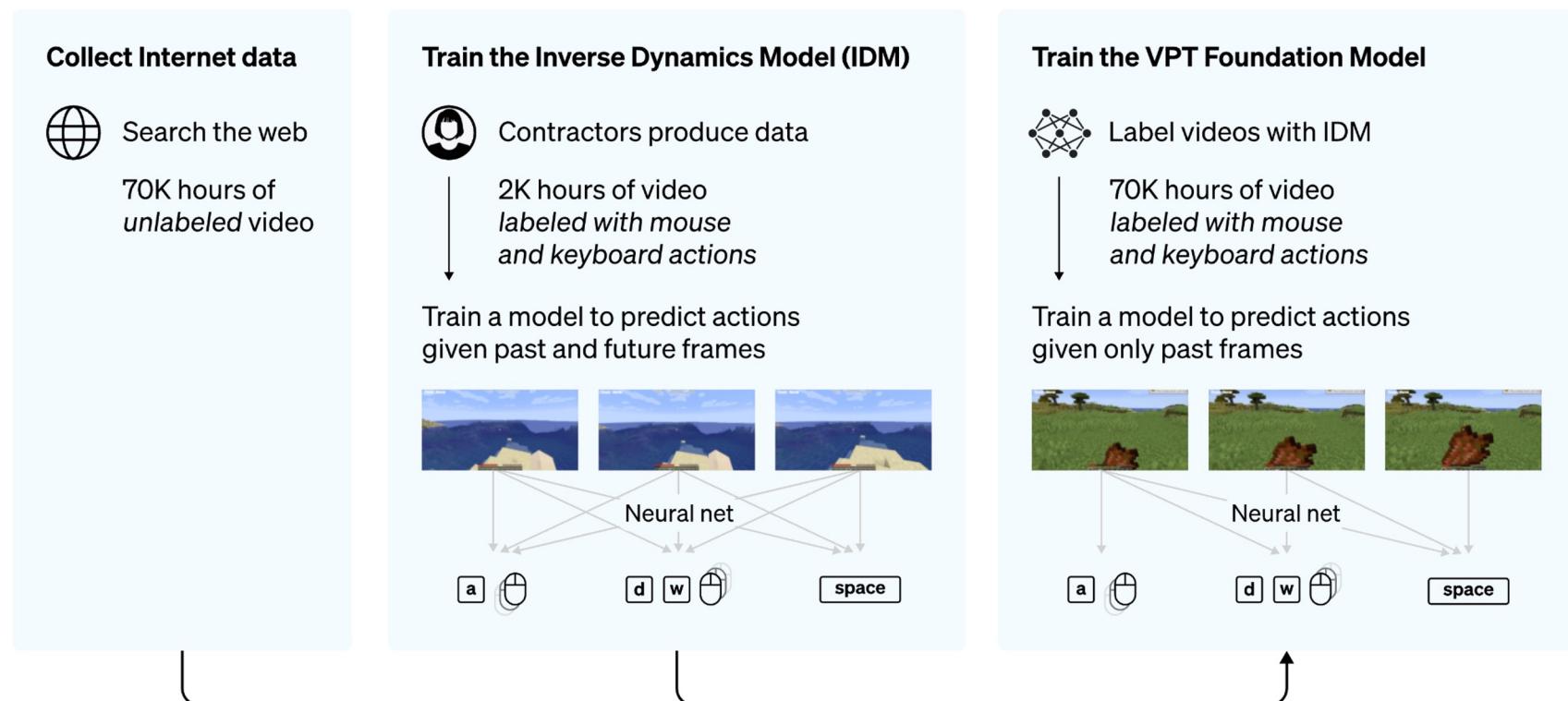
	LP Search	LP Choice	Human Demo	Use Reward
Rule				
IL w/o LP Choice	✓		✓	
IL w/o LP Search		✓	✓	
IL	✓	✓	✓	
RL	✓			✓
RL (RNN)				✓
IL+RL	✓	✓	✓	✓



Click and play

# Interactive space: Sandbox

- Video pre-training of MineCraft
  - Sandbox like Minecraft is a good interactive space





# Interactive space: Sandbox

- Video pre-training of MineCraft
  - Define discrete actions in the interactive space

Action	Human action	Description
forward	W key	Move forward.
back	S key	Move backward.
left	A key	Strafe left.
right	D key	Strafe right.
jump	space key	Jump.
inventory	E key	Open or close inventory and the 2x2 crafting grid.
sneak	shift key	Move carefully in current direction of motion. In the GUI it acts as a modifier key: when used with attack it moves item from/to the inventory to/from the hot-bar, and when used with craft it crafts the maximum number of items possible instead of just 1.
sprint	ctrl key	Move fast in the current direction of motion.



# Interactive space: Sandbox

- Cost
  - Use behavior model to annotate unlabeled 70 hours video
  - Reduce the cost: 1,400,000 \$ -> 130,000 \$
- Annotation trick
  - At first, casually playing MineCraft
  - Play specific tasks (Equip Capabilities)

***Collect as many units of wood as possible, using only wooden or stone tools***

***Start a new world every 30 minutes of game play***

***Build a basic house in 10 minutes using only dirt, wood, sand, and either wooden or stone tools***

***Starting from a new world and an empty inventory, find resources and craft a diamond pickaxe in 20 minutes***



# Interactive space: Sandbox

- Results
  - VPT accomplishes tasks impossible to learn with RL alone, such as **crafting planks** and **crafting tables** (tasks requiring a human proficient of ~970 consecutive actions)





# Interactive space: Sandbox

- Results
  - An example for killing a cow





# Challenges & limitations

- Challenges of the current framework
  - Building **interactive space** is time-consuming
  - **Labeling** is expensive and labor-intensive
  - The goal must be clear and simple
  - Only **discrete** actions and states are supported
  - A clean interactive space is required

Goal	Space	Actions
Why do people cry?		
What's the most common occupations of American's first lady		
Model an awesome ironman;		



# Thanks for listening

THUNLP