

# Similarity Analysis of Contextual Word Representation Models

Wu+ ACL 2020

# Motivation

- There are many contextual word representations proposed based on different architectures.
  - ElMo, BERT, XLNet, GPT, Albert, RoBERTa, ...
- How are these different empirically?
- This paper proposes various methods to measure the similarity between the models trained using different architectures.
- This is different from probing where models are compared based on their performance on different downstream (or probing) tasks
  - We need linguistic annotations for probing
  - Comparisons between models are indirect (via probing accuracy)
- Research questions addressed in this paper
  - Do different models behave similarly on the same inputs?
  - Which design choices determine whether models behave similarly or differently?
  - Are certain model components more similar than others across architectures?
  - Is information in a model localised (i.e. one computing element for each represented entity)

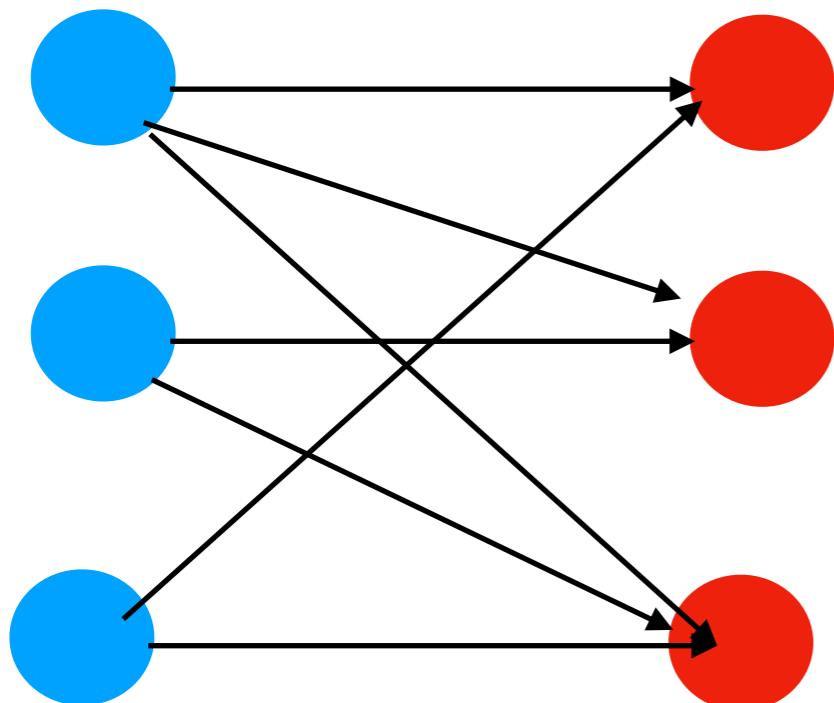
# Findings

- Different architectures have rather similar representations, but different individual neurones.
- Lower layers are more similar than higher layers across architectures.
- Higher layers have more localised representations than lower layers.
- Higher layers are more affected by fine-tuning on downstream tasks than lower layers.
- Fine-tuning affects the localisation of information, causing high layers to be less localised.

# Similarity Measures

- Collection of M models  $\{f^{(m)}\}_{m=1}^M$
- word representations  $\mathbf{h}_l^{(m)}$  by the m-th model in its l-th layer
- attention weights  $\alpha_l^{(m)}$
- Neurones are indexed by k,  $\mathbf{h}_l^{(m)}[k]$  and  $\alpha_l^{(m)}[k]$
- In the case of word embeddings [k] denotes the k-th dimension and for attention weights [k] will be the value of the k-th attention head (possibly a matrix)
- Various similarity measures are proposed between neurones,  $\text{sim}(\mathbf{h}_l^{(m)}, \mathbf{h}_{l'}^{(m')})$  and between attentions  $\text{sim}(\alpha_l^{(m)}, \alpha_{l'}^{(m')})$

# Neurones vs. Dimensions



hidden layer  
dimensions correspond to neurones

We are not interested in the edges (weights)  
in this analysis

# Neurone-level Similarity

Compare one neurone from a layer of one model against another neurone from a layer in a different model.

For a given neuron  $\mathbf{h}_l^{(m)}[k]$ , we define

$$\text{neuronsim}(\mathbf{h}_l^{(m)}[k], \mathbf{h}_{l'}^{(m')}) = \max_{k'} |\rho(\mathbf{h}_{l'}^{(m')}[k'], \mathbf{h}_l^{(m)}[k])|$$

as the maximum correlation between it and another neuron in some layer (Bau et al., 2019). Here  $\rho$  is the Pearson correlation. This naturally gives rise to an aggregate measure at the layer level:

$$\text{neuronsim}(\mathbf{h}_l^{(m)}, \mathbf{h}_{l'}^{(m')}) = \frac{1}{N_m} \sum_k \text{neuronsim}(\mathbf{h}_l^{(m)}[k], \mathbf{h}_{l'}^{(m')})$$

# Mixed neurone-representation similarity

Information of one neurone can be distributed across neurones in another layer.  
Do least square regression and find the minimum regression score.

We define

$$\text{mixedsim}(\mathbf{h}_l^{(m)}[k], \mathbf{h}_{l'}^{(m')}) := \\ \text{lstsq}(\mathbf{h}_{l'}^{(m')}, \mathbf{h}_l^{(m)}[k]) \cdot r$$

where  $\cdot r$  is the r-value associated with the regression, the norm of the prediction divided by the norm of the regressand. As before, this is extended to the layer level:

$$\text{mixedsim}(\mathbf{h}_l^{(m)}, \mathbf{h}_{l'}^{(m')}) = \\ \frac{1}{N_m} \sum_k \text{mixedsim}(\mathbf{h}_l^{(m)}[k], \mathbf{h}_{l'}^{(m')})$$

# Representation-level Similarity

- A representation-level measure find correlations between a full model (or layer) simultaneously.
  - Canonical Correlation Analysis (CCA) [Raghuram et al. 2017]
  - Projection weighted CCA [Kornblith et al. 2019]

**SVCCA** Given two layers

$$\mathbf{X}, \mathbf{Y} = \mathbf{Z}\mathbf{h}_{l_x}^{(m_x)}, \mathbf{Z}\mathbf{h}_{l_y}^{(m_y)}$$

we compute the truncated principal components

$$\mathbf{X}', \mathbf{Y}' = \mathbf{U}_x[:, :l_x], \mathbf{U}_y[:, :l_y]$$

where  $\mathbf{U}_x$  are the left singular vectors of  $\mathbf{X}$ , and  $l_x$  is the index required to account for 99% of the variance.  $\mathbf{U}_y$  and  $l_y$  are defined analogously. The SVCCA correlations,  $\rho_{SVCCA}$ , are defined as:

$$\mathbf{u}, \rho_{SVCCA}, \mathbf{v} = \text{SVD}(\mathbf{X}'^T \mathbf{Y}')$$

The SVCCA similarity,  $\text{svsim}(\mathbf{h}_{l_x}^{(m_x)}, \mathbf{h}_{l_y}^{(m_y)})$ , is the mean of  $\rho_{SVCCA}$ .

**PWCCA** Identical to SVCCA, except the computation of similarity is a weighted mean. Using the same notation as above, we define canonical vectors,

$$\mathbf{H}_X := \mathbf{X}'\mathbf{u}$$

$$\mathbf{H}_Y := \mathbf{Y}'\mathbf{v}$$

We define alignments

$$\mathbf{A}_X := \text{abs}(\mathbf{H}_X^T \mathbf{X})$$

$$\mathbf{A}_Y := \text{abs}(\mathbf{H}_Y^T \mathbf{Y})$$

where  $\text{abs}$  is the element-wise absolute value. The weights are

$$\alpha_x := \text{weights}(\mathbf{A}_X \mathbf{1}), \quad \alpha_y := \text{weights}(\mathbf{A}_Y \mathbf{1})$$

where  $\mathbf{1}$  is the column vector of all ones, and  $\text{weights}$  normalizes a vector to sum to 1. The PWCCA similarity is

$$\text{pwsim}(\mathbf{h}_{l_x}^{(m_x)}, \mathbf{h}_{l_y}^{(m_y)}) := \alpha_x^T \rho_{SVCCA}$$

$$\text{pwsim}(\mathbf{h}_{l_y}^{(m_y)}, \mathbf{h}_{l_x}^{(m_x)}) := \alpha_y^T \rho_{SVCCA}$$

It is asymmetric.

# Attention-level Similarity

- Find the correlation between attention heads
  - Norm between the two attention matrices  $\alpha_l^{(m)}[k], \alpha_{l'}^{(m')}[k']$
  - Their Pearson correlation
  - Their Jensen-Shanon divergence

We define

$$\text{attnsim}(\alpha_l^{(m)}[k], \alpha_{l'}^{(m')}) = \max_{k'} [\text{Sim}(\alpha_{l'}^{(m')}[k'], \alpha_l^{(m)}[k])]$$

We consider three such values of Sim.

- Matrix norm: for each sentence  $s_i$ , compute the Frobenius norm  $\|\alpha_{l'}^{(m')}[h'](s_i) - \alpha_l^{(m)}[h](s_i)\|$ . Then average over sentences in the corpus.
- Pearson correlation: for every word  $x_i$ , compare the attention distributions the two heads

induce from  $x_i$  to all words under Pearson correlation:  $\rho(\alpha_{l',i}^{(m')}[h'], \alpha_{l,i}^{(m)}[h])$ . Then average over words in the corpus.

- Jensen–Shannon divergence: for every word  $x_i$ , compare the attention distributions under Jensen–Shannon divergence:  $\frac{1}{2} \text{KL}(\alpha_{l',i}^{(m')}[h'] \parallel \beta) + \frac{1}{2} \text{KL}(\alpha_{l,i}^{(m)}[h] \parallel \beta)$ , where  $\text{KL}$  is the KL-divergence and  $\beta$  is the average of the two attention distributions. Then average of words in the corpus.

# Distributed Attention-level Similarity

- Compare entire attention heads in two layers by concatenating all weights from all heads in one layer to get an attention representation.

# Experimental Setup

- Models
  - ELMo variants: original, bidirectional RNN with 2 hidden layers, 4 layers and Transformer-equivalent variant
  - GPT variants: GPT, GPT-2
  - BERT: base, large
  - XLNet: base, large

# Data

- Run on Penn Treebank development and obtain word representations
- Aggregated sub-word representations by taking the representation fr the last sub-word
- Sub-word attentions are computing by summing up attention to sub-words and averaging attention from sub-words (guarantees that attention from each word sums to one)

# Similarity of Pre-trained Models

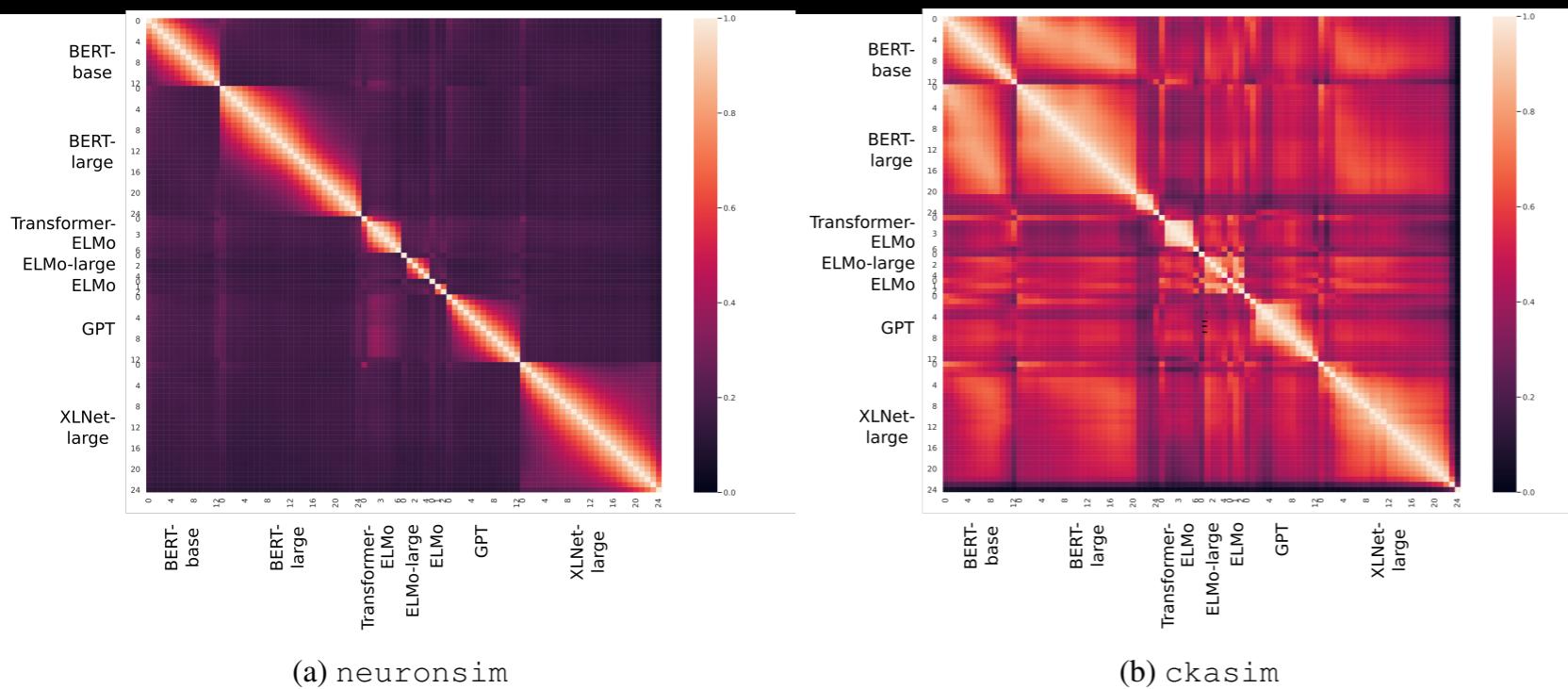
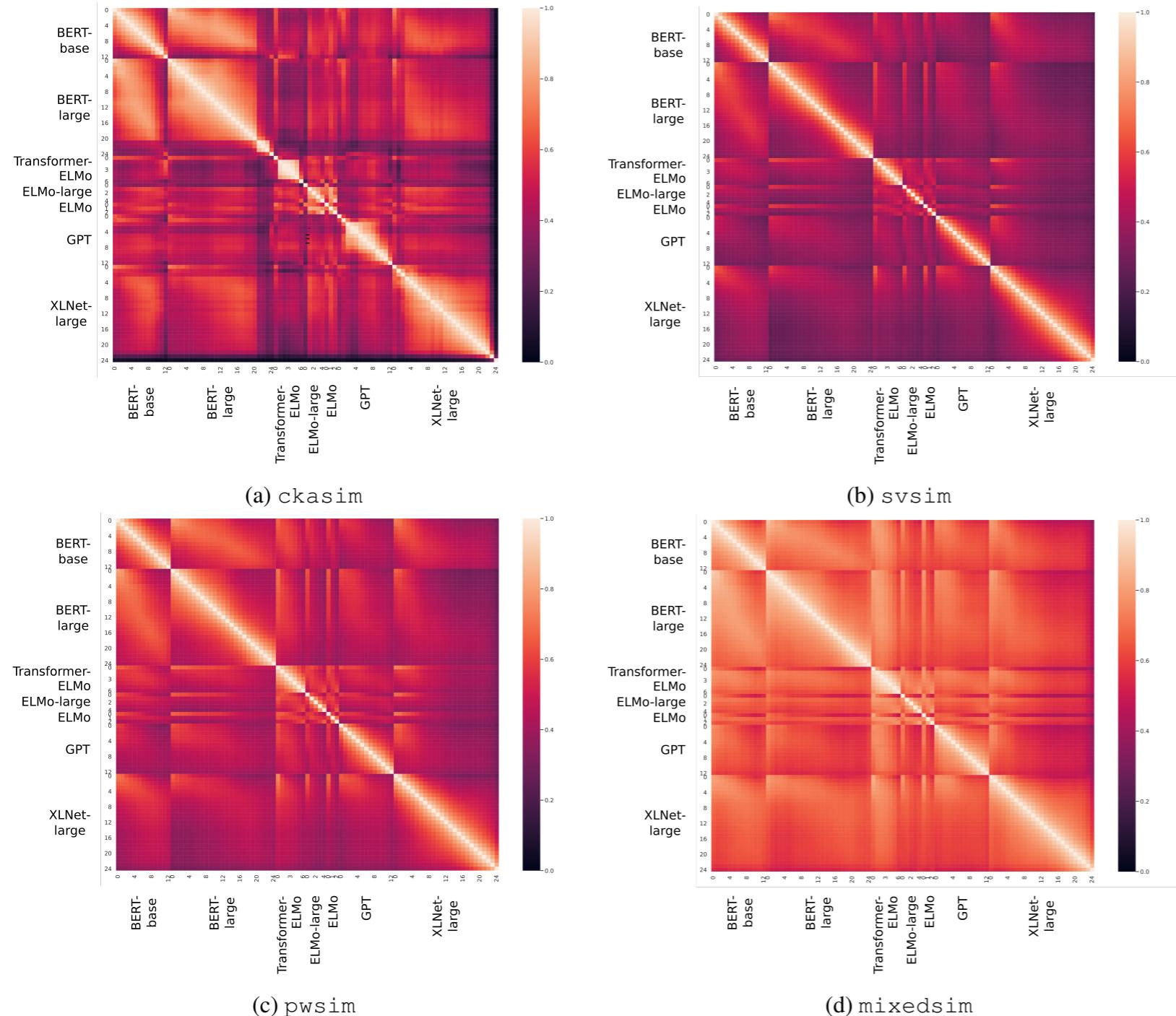


Figure 1: Similarity heatmaps of layers in various models under neuron- and representation-level similarities.

- **neurosim is block diagonal — within a model, different layers have similar individual neurones, but across models, neurones are very different.**
- **ckasim has high off-diagonal values as well — different models generate similar representations**
- Transformer-ELMo shares ELMo's bidirectional objective function but with Transformers rather than RNNs. It is similar to both ELMo and GPT, more so than BERT or XLNet.

**ckasim:** representation-level similarity based on linear centred kernel alignment

# mixedsim

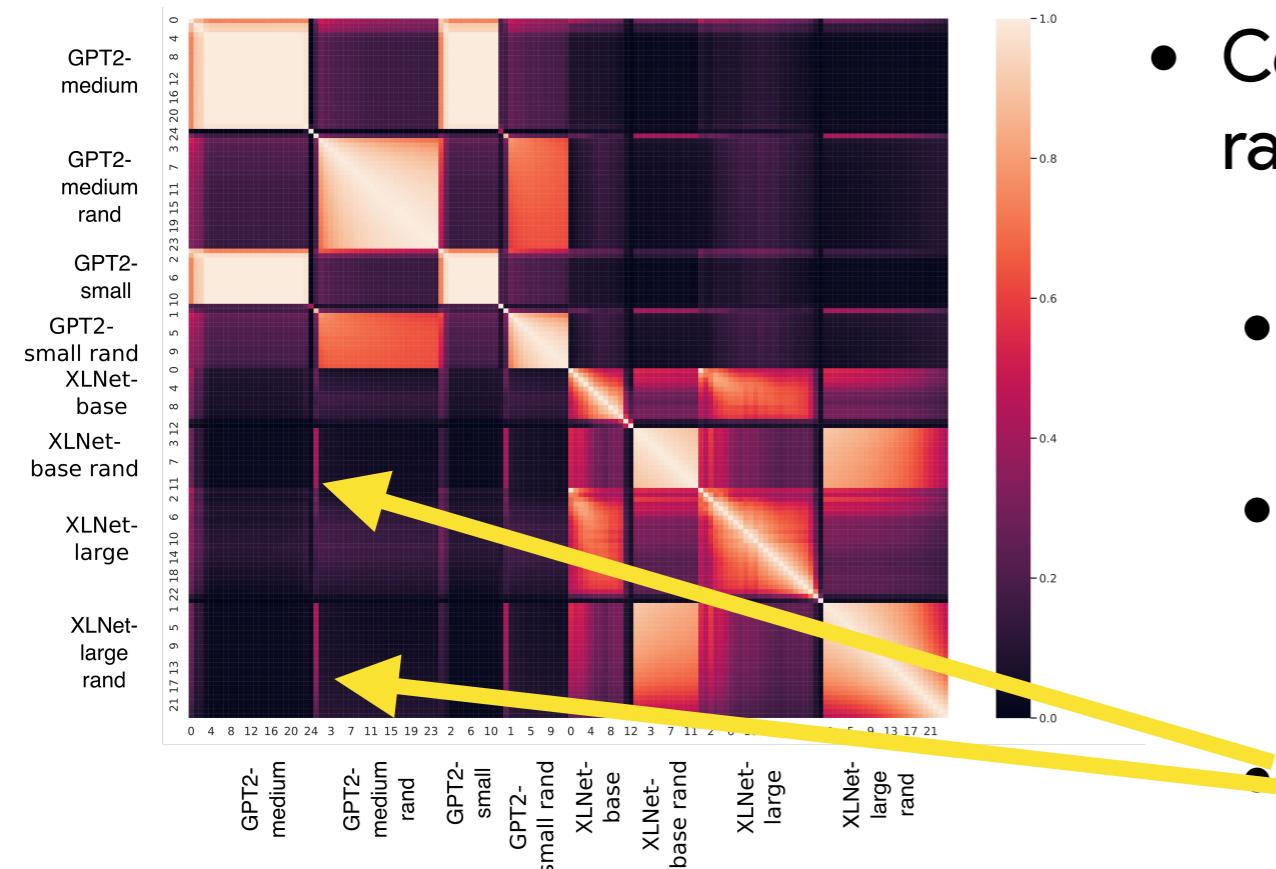


- Individual neurones in one model may be well represented by a linear combination of neurones in another layer
- Models within the same family are more similar

Figure 8: Similarity heatmaps of layers in various models under different representation-level similarity measures.

# Why?

- Why are models within the same family more similar to each other?
  - Could be because models within the same family are often trained on the same data, but cross-family models are trained on different data.

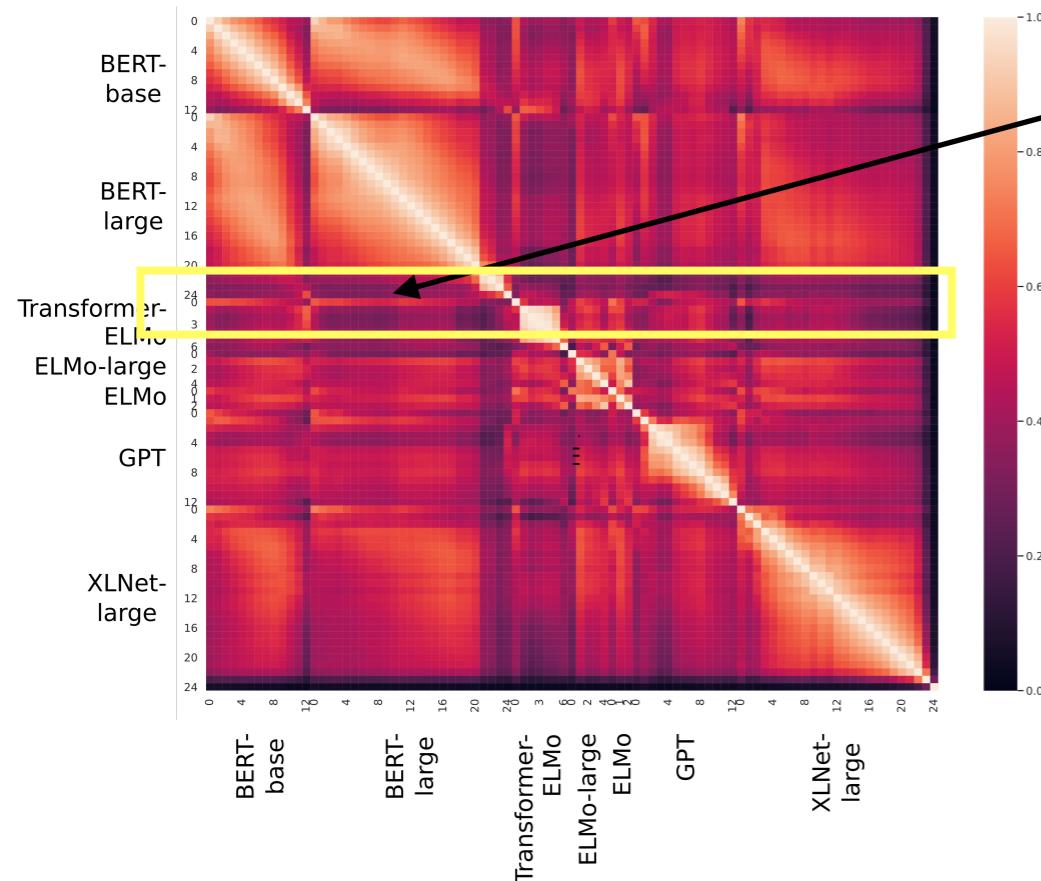


- Comparing layers from pre-trained and randomly initialised models
  - Still high intra-model similarities
  - Random vs. pre-trained models are not very similar (as expected)
- Vertical bands in lower triangle

Figure 2: ckasim similarity heatmap of layers in base and random models.

- Lower layers converge faster, leaving them closer to their initial random state.

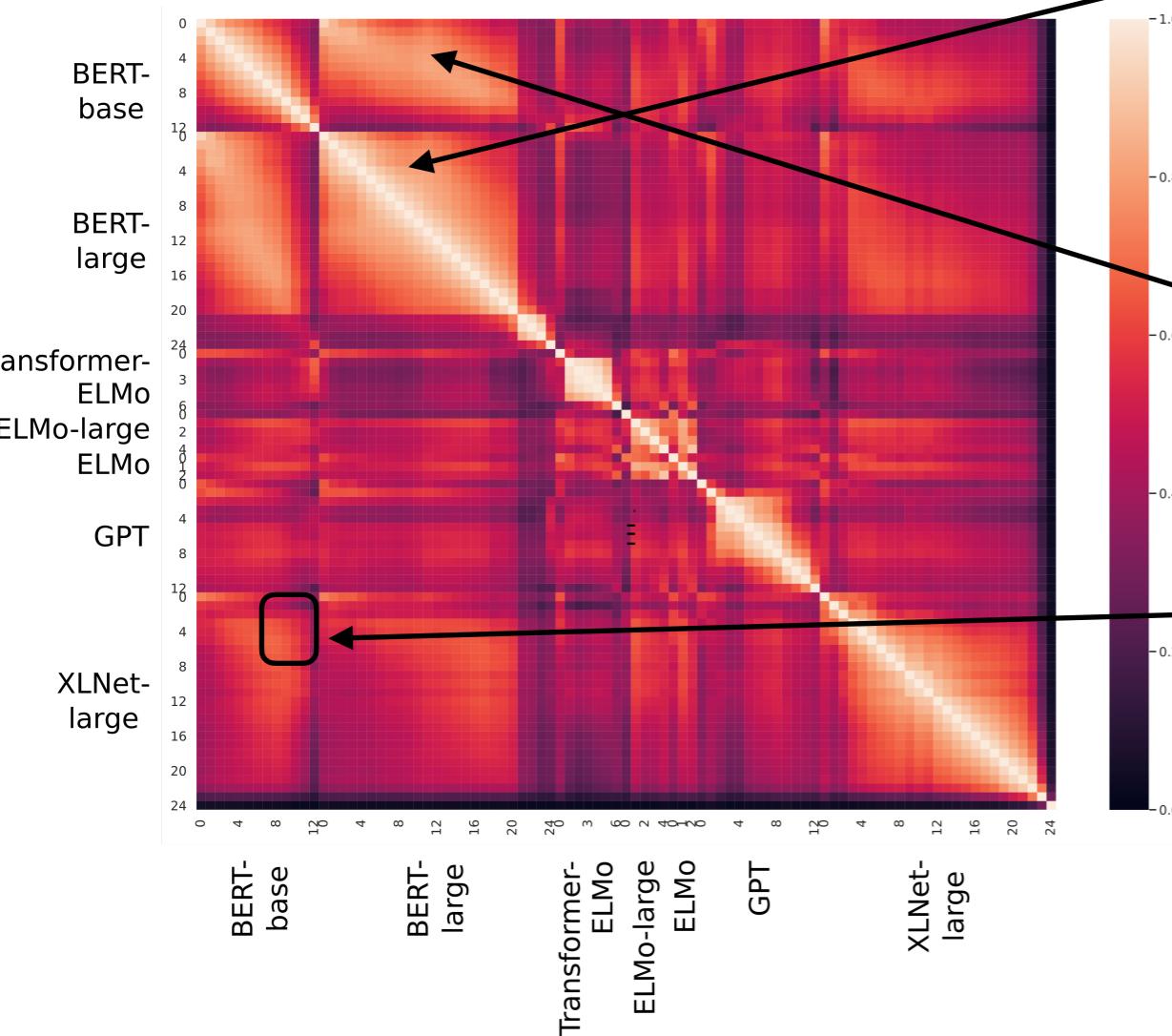
# Lower layers



(b) ckasim

- Lower layers are more similar across architectures
- Lower layers are closer to the input, which is always the same words

# Adjacent layers are more similar



(b) ckasim

- There are “thick” bright diagonals, meaning that adjacent layers are more similar.
- BERT-base is similar to BERT-large as well
- Lower layers of XLNet are more similar to higher layers of BERT
- BERT requires mostly features from higher layers to achieve SoTA, when XLNet need lower and middle layers (Liu+19)

# Higher layers are more localised than lower ones

- Average the neuronsim (local) and svssim (distributed) similarities between a layer and all the other layers. Subtract avg. neuronsim from avg. svssim to compute a “localisation score”

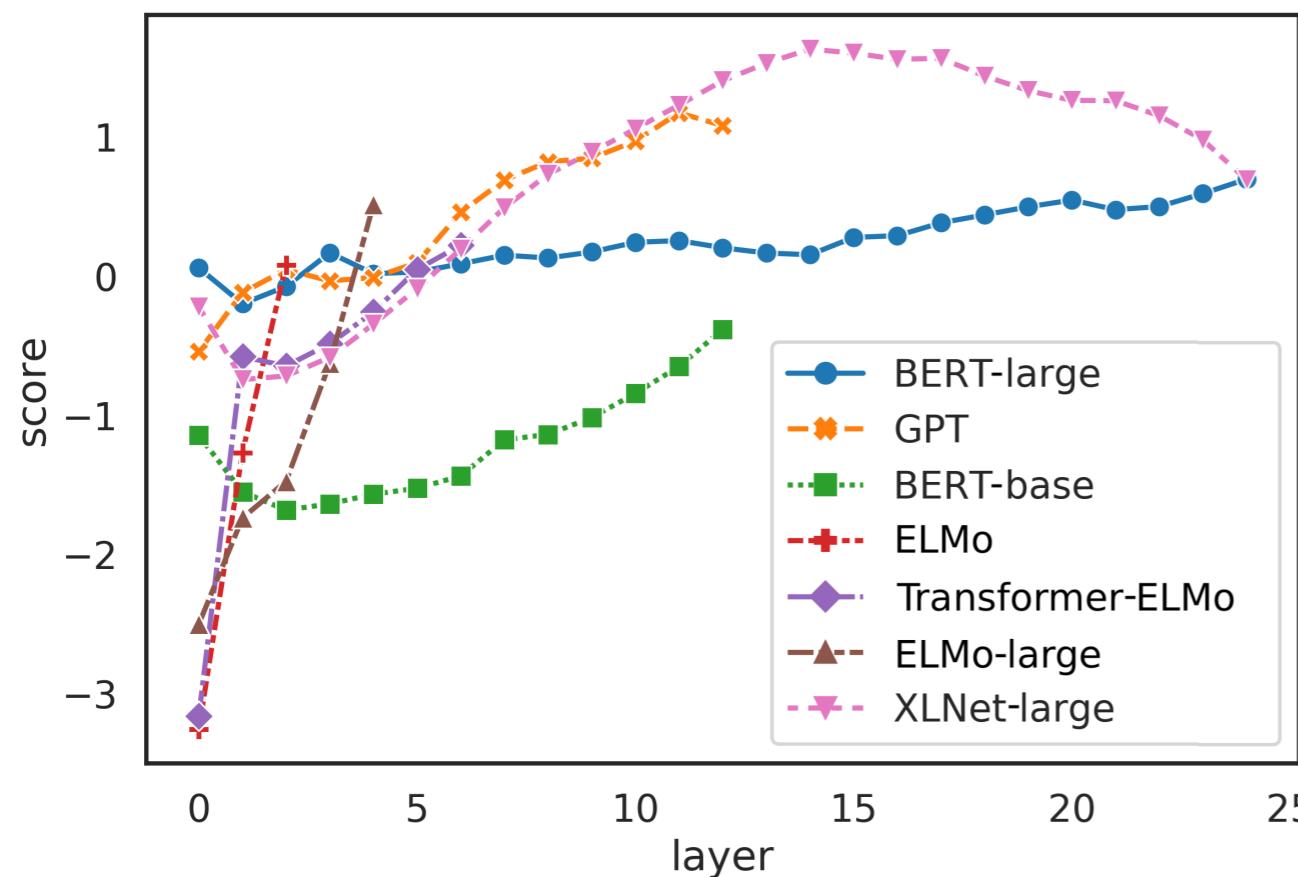
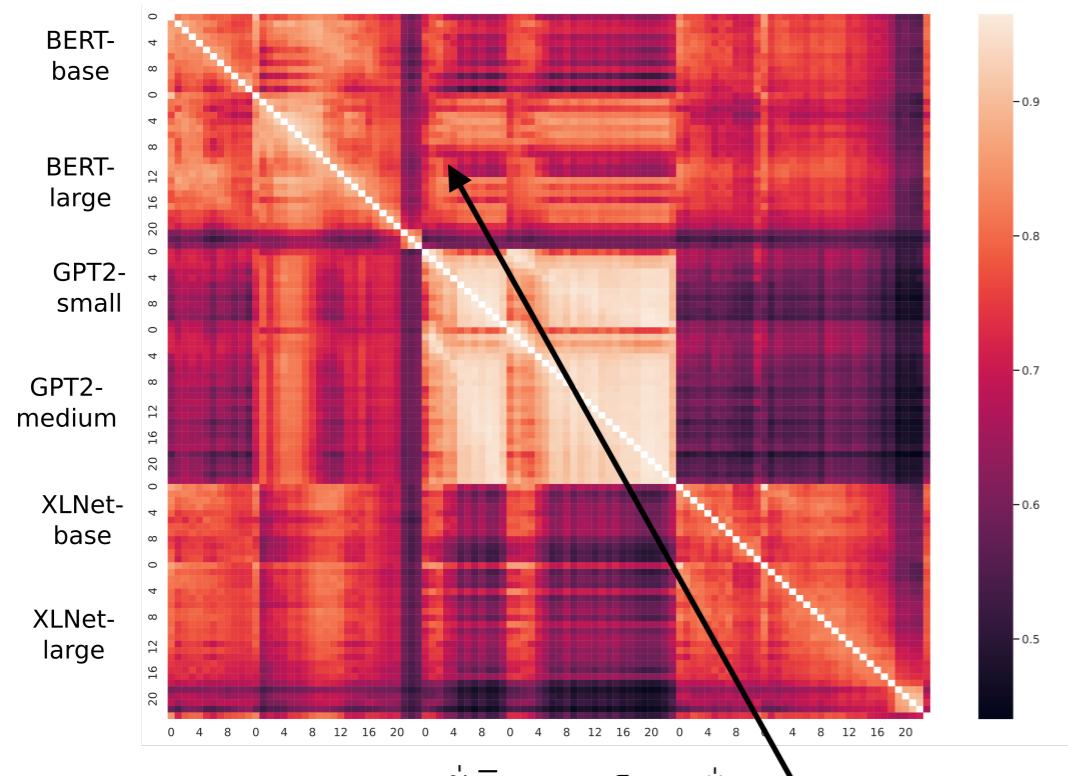


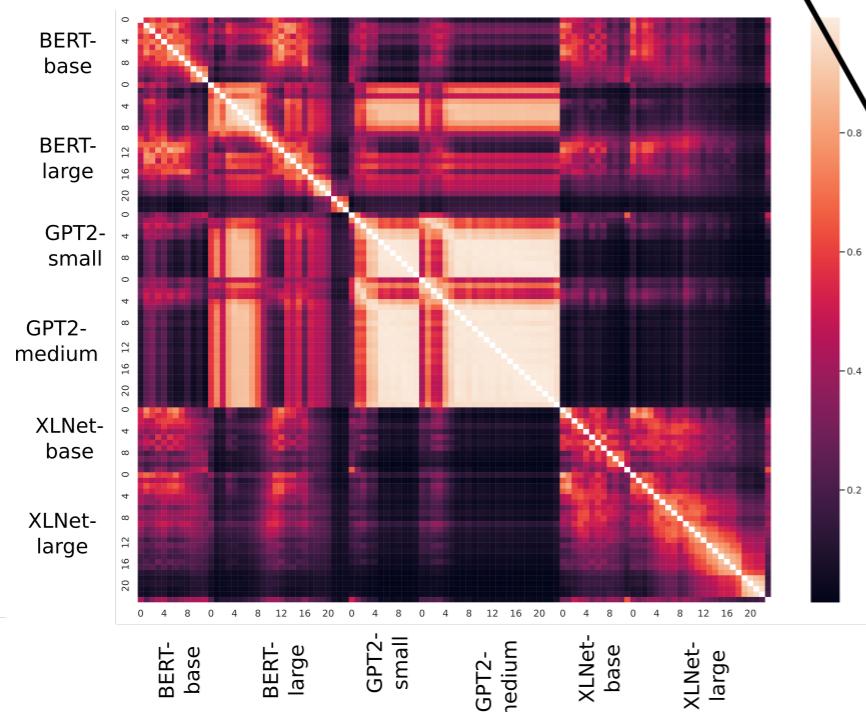
Figure 3: Localization score of various model layers.

- In all models, localisation score increases with layers
- Upper layers are more context-specific (Ethyarajah+19) — Same word's embeddings are less similar in different contexts
- localised  $\approx$  context-specific
- XLNet’s top layers are less context-specific

# Attention Level



(a) Jensen–Shannon



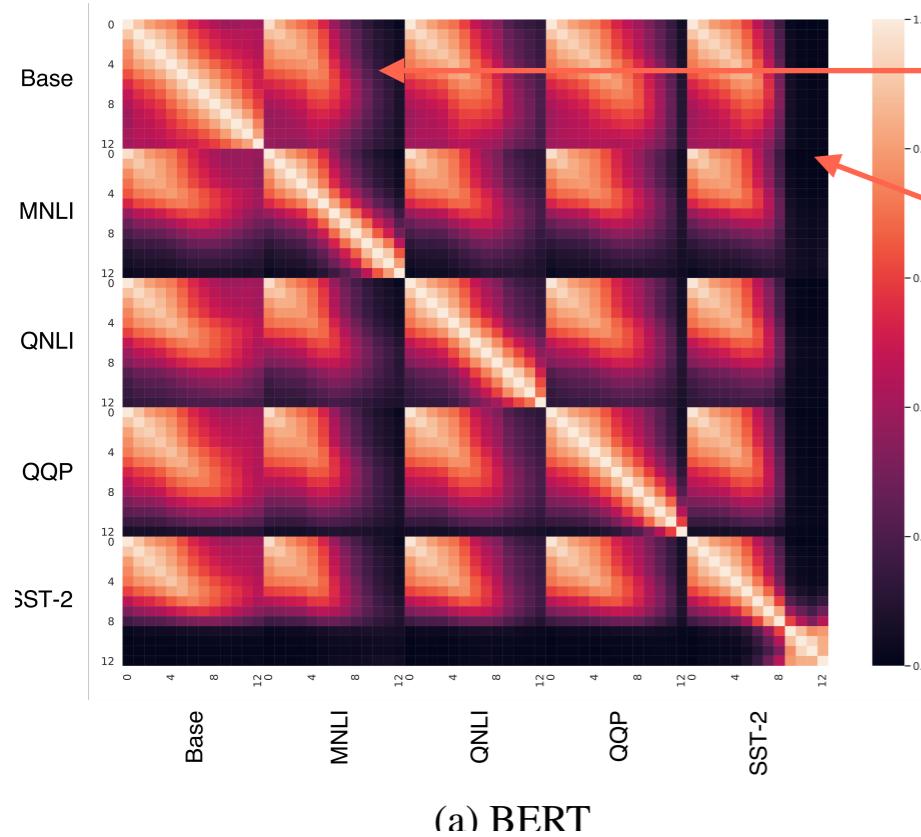
(b) ckasim

- Layers within the same model or model family are highly similar
- GPT2 layers are all similar to each other
- JS shows higher non-diagonal similarities than ckasim
- Difficult to identify patterns within a given model family (attention is difficult to interpret)
- GPT2 layers are similar to the bottom layers of BERT-large

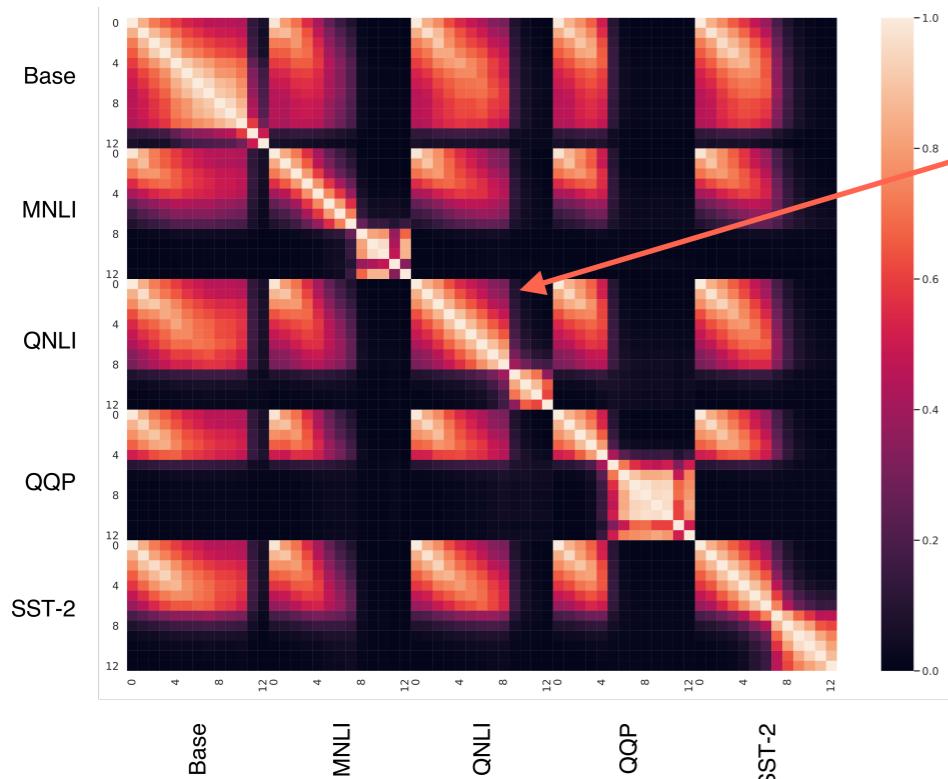
# Similarity of Fine-tuned Models

- How does fine-tuning on downstream tasks affect model similarity?
  - Tasks:
    - MNLI (multilingual NLI)
    - QNLI (predict whether a sentence contains the answer to a question)
    - QQP: Classify whether two (Quora) questions are semantically equivalent (duplicate detection)
    - SST-2: binary sentiment classification on Stanford sentiment treebank

# Similarity of Fine-tuned Models

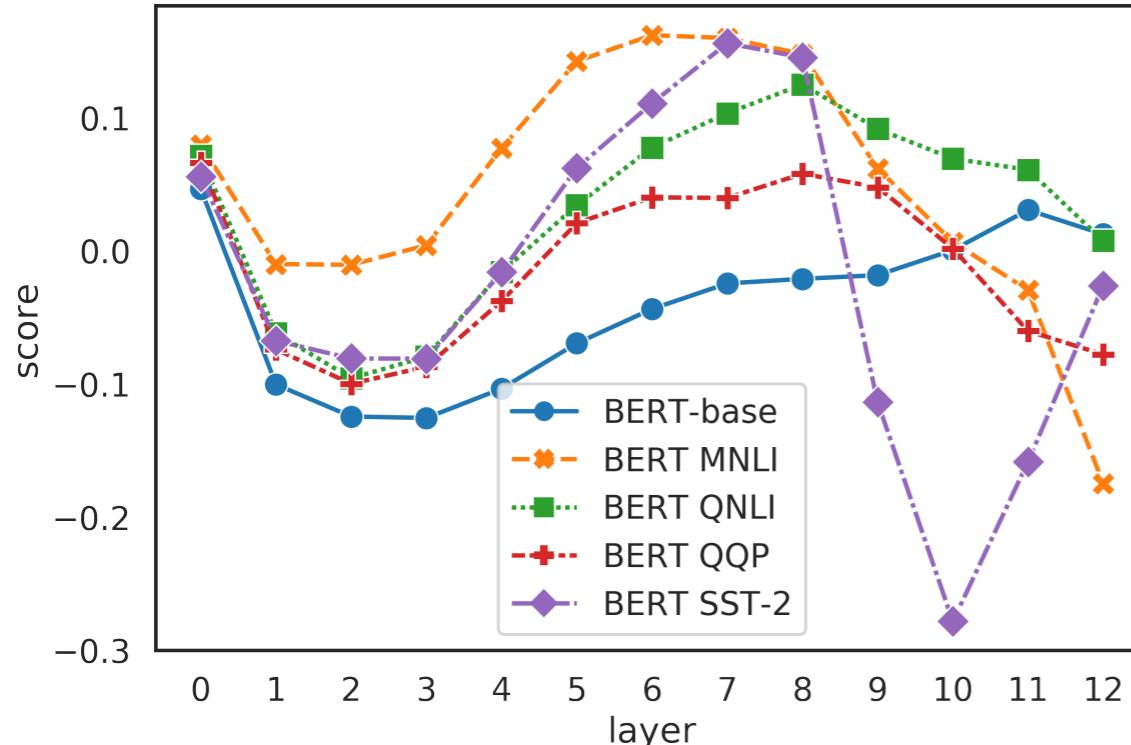


- Top layers are more affected (becomes less similar) after fine-tuning
- In BERT, top layers of SST-2 fine-tuned model are affected more
  - SST-2 is sentence classification whereas other tasks are sentence-pair classification

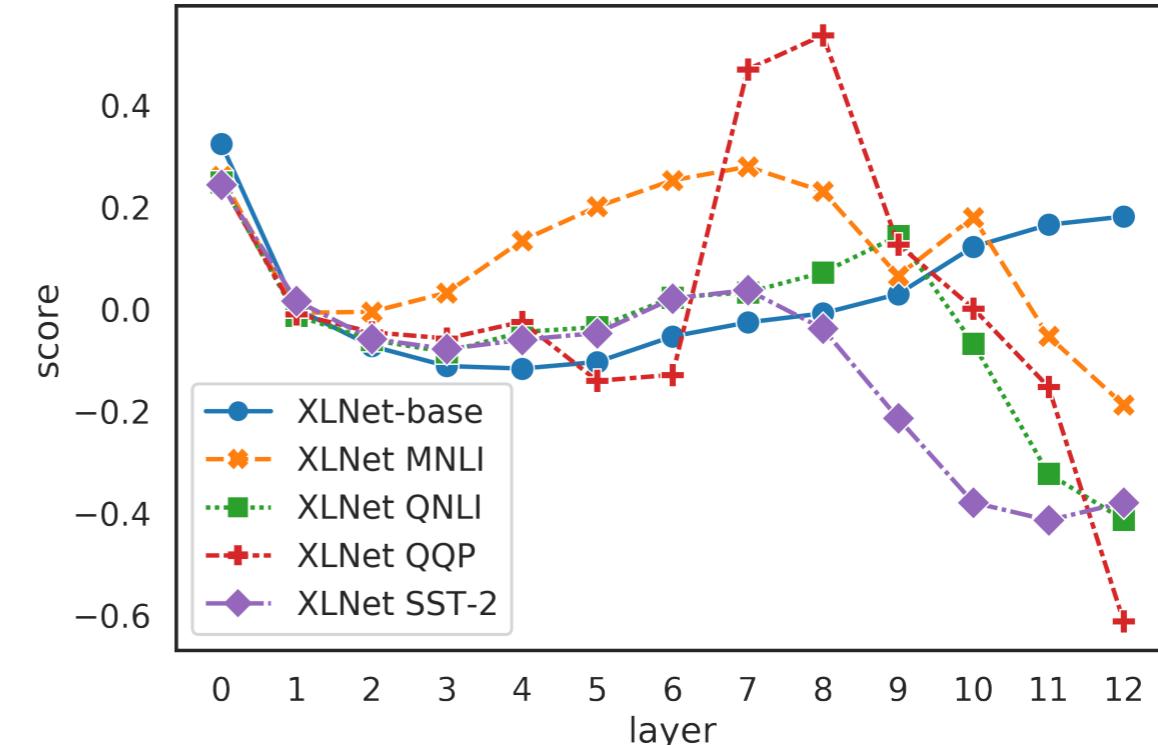


- In XLNet, fine-tuning on any task leads to top layers being very different from all layers of models fine-tuned on other tasks
- Difficult to multitask with XLNet compared to BERT?

# Fine-tuning affects localisation



(a) BERT



(b) XLNet

- Compared to the pre-trained models, the fine-tuned ones decrease in localisation at the top layers.
  - To learn high-level tasks, we might require multiple specialised neurones.

# Future directions

- Combine probing with similarity analysis (this paper) to identify which linguistic properties are captured in the model components that are similar to one another.
- Relationship between localisation of information and learnability of particular properties
- Context specificity vs. localisation connection
- Attention similarity vs. Interpretability/Explainability
- Even if two models have similar layers, the outputs could be very different because how those layers are connected (recurrent vs. transformer)