

Deconfounded Representation Similarity for Comparison of Neural Networks

Tianyu Cui¹, Yogesh Kumar¹, Pekka Marttinen¹, Samuel Kaski^{1,2}

¹Aalto University ²The University of Manchester

BACKGROUND:

Representation similarities between NNs:

For input $X \in \mathbb{R}^{n \times p}$ and the corresponding m_1 th and m_2 th layer representations of NNs f_1 and f_2 , $X_{f_1}^{m_1}$ and $X_{f_2}^{m_2}$, we

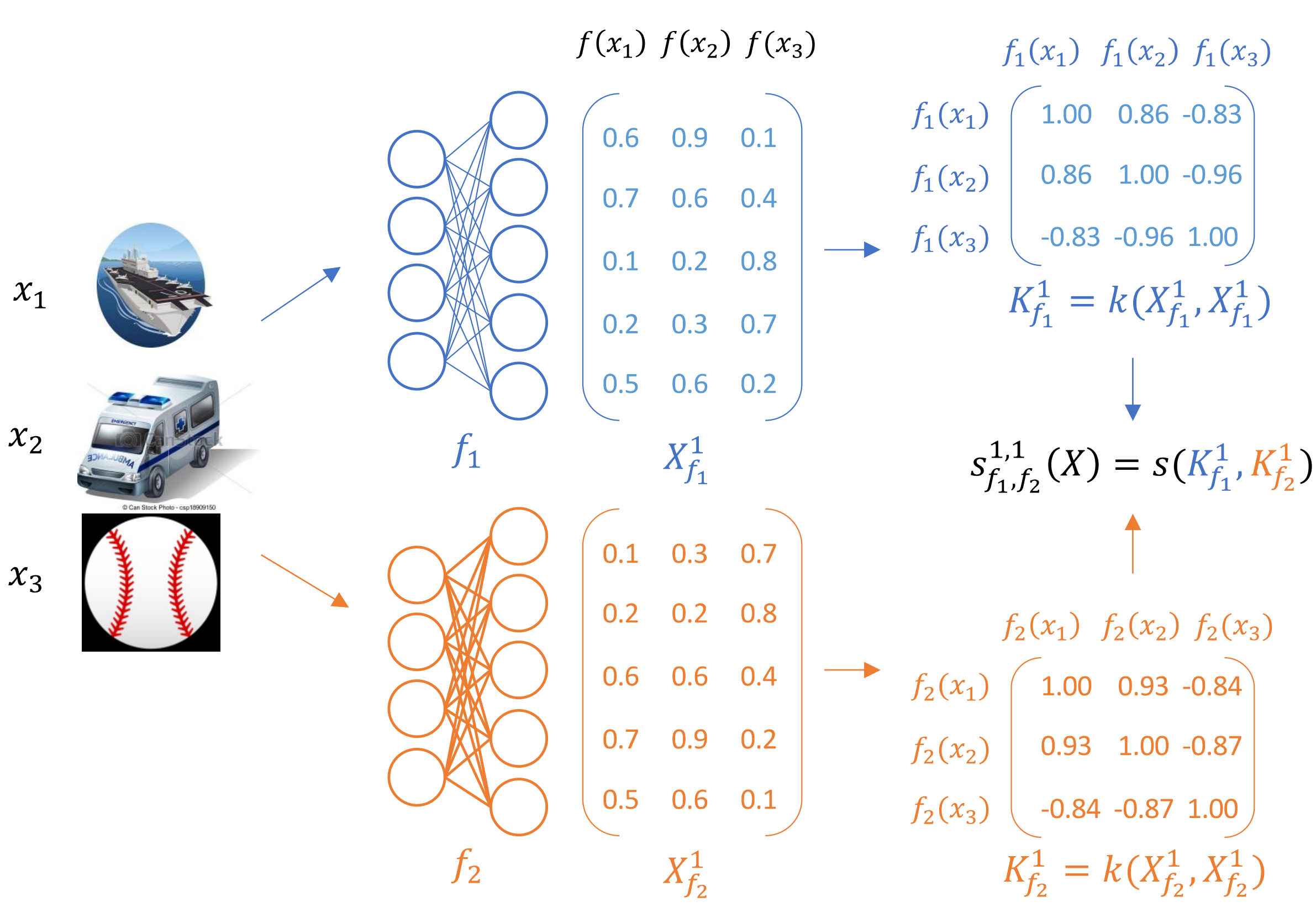
1. Compute the inter-example representation similarity matrices (RSMs) of $X_{f_1}^{m_1}$ and $X_{f_2}^{m_2}$:

$$K_{f_1}^{m_1} = k(X_{f_1}^{m_1}, X_{f_1}^{m_1}), K_{f_2}^{m_2} = k(X_{f_2}^{m_2}, X_{f_2}^{m_2});$$

2. Compute the similarity between two RSMs:

$$s_{f_1, f_2}^{m_1, m_2} = s(K_{f_1}^{m_1}, K_{f_2}^{m_2}).$$

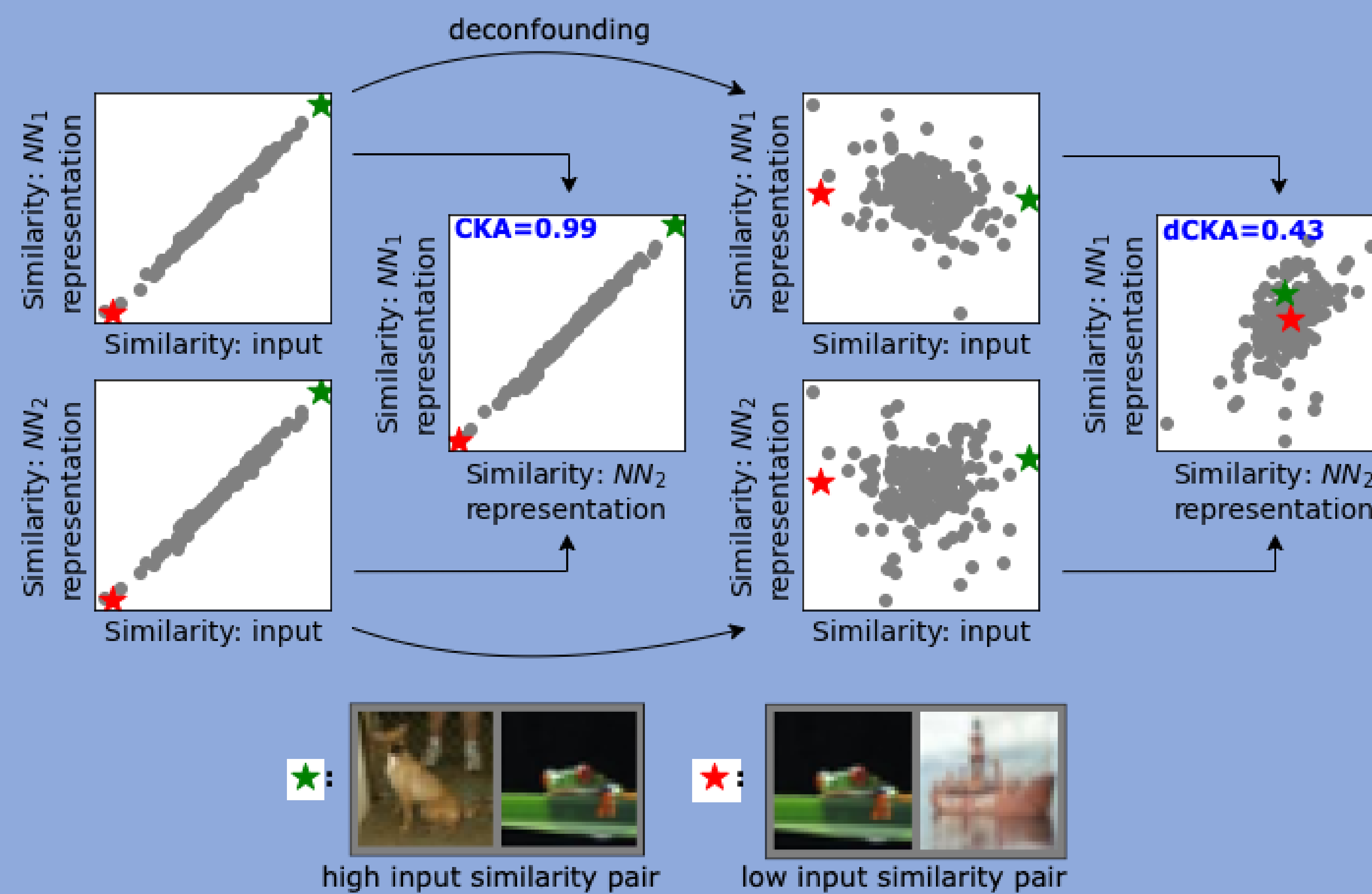
This covers most of methods, such as CKA¹.



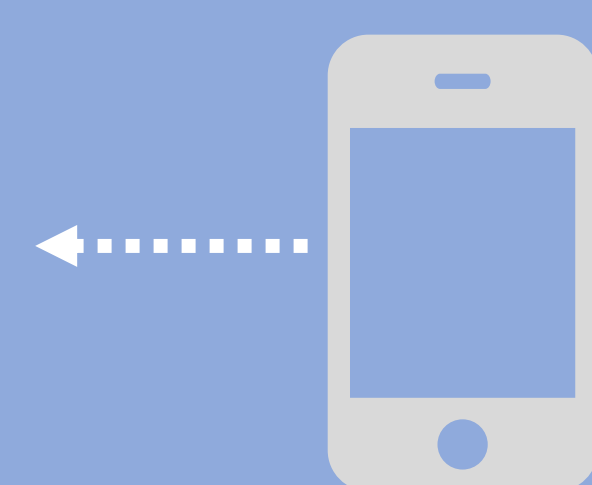
How to better compare representations learned by neural networks?

Current measures for comparing the similarity of representations between NNs may not well reflect their functional similarity due to the confounding effect of similarities of data items in the input space.

Comparing Representations of Random ResNets (NN_1, NN_2)



Counterintuitively, a conventional method (CKA¹) gives completely random NNs a similarity close to one. Our dCKA gives a smaller similarity.



Take a picture to access the paper and code!

METHODS:

Deconfounded representation similarities:

Just regress out the confounder from RSMs:

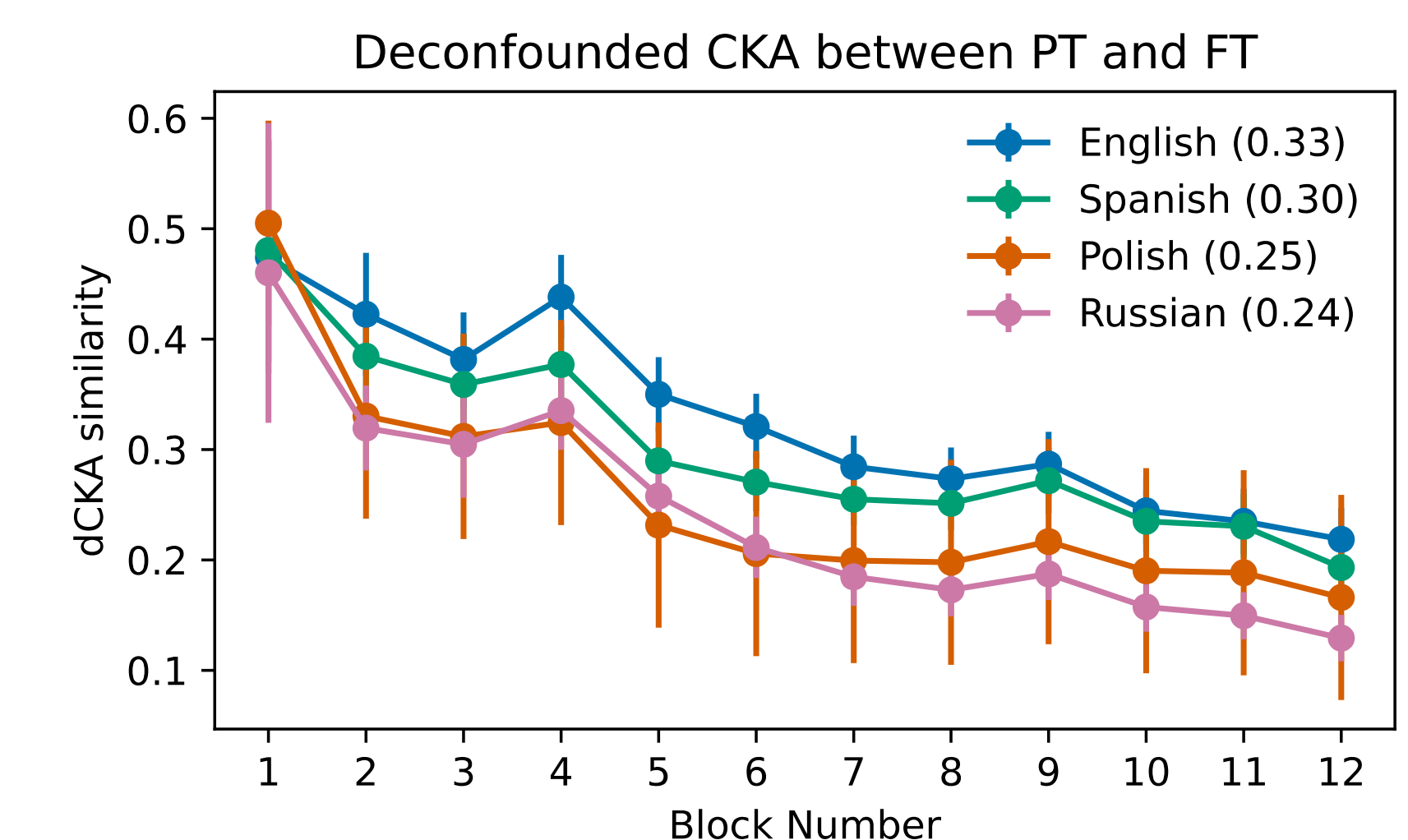
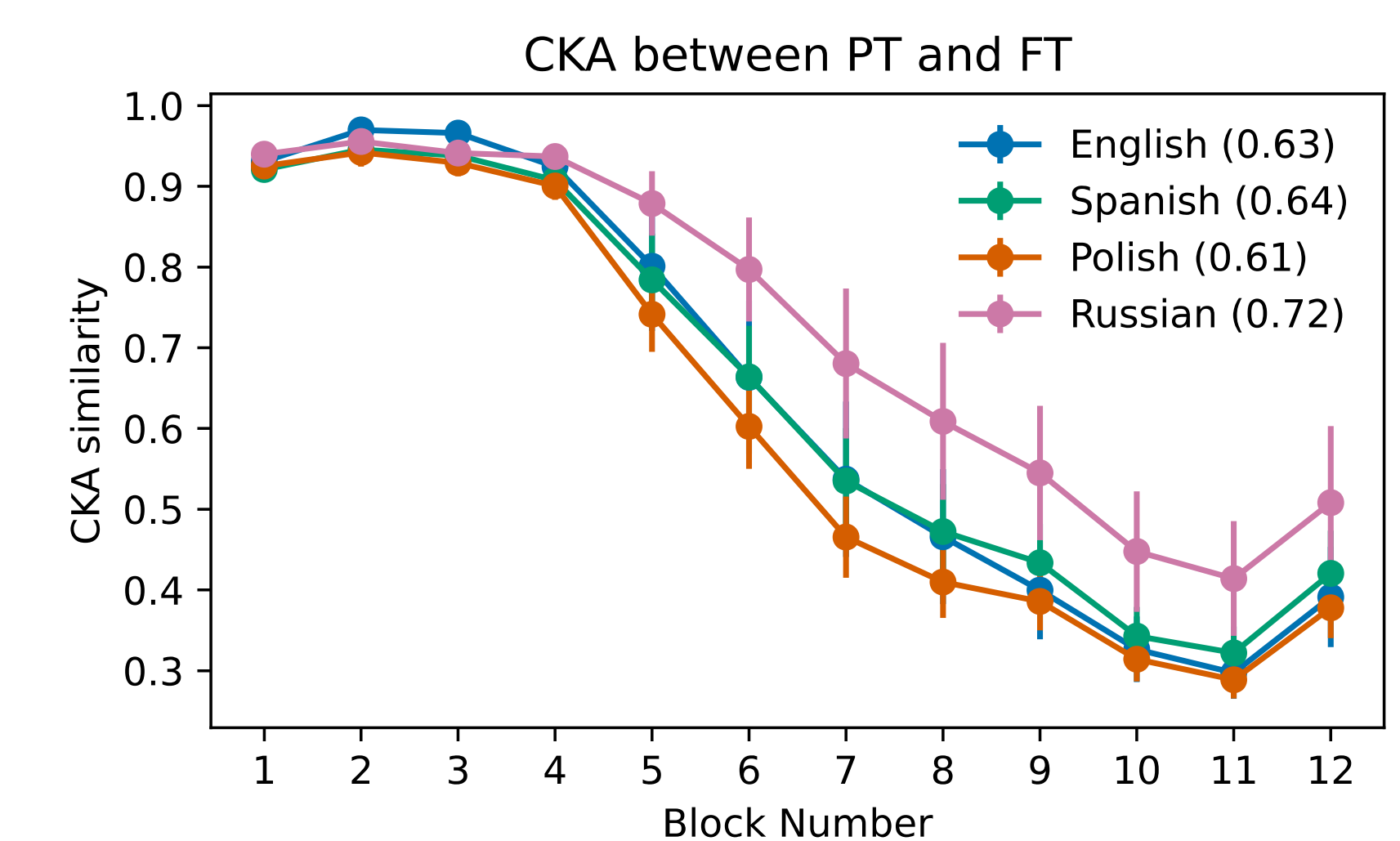
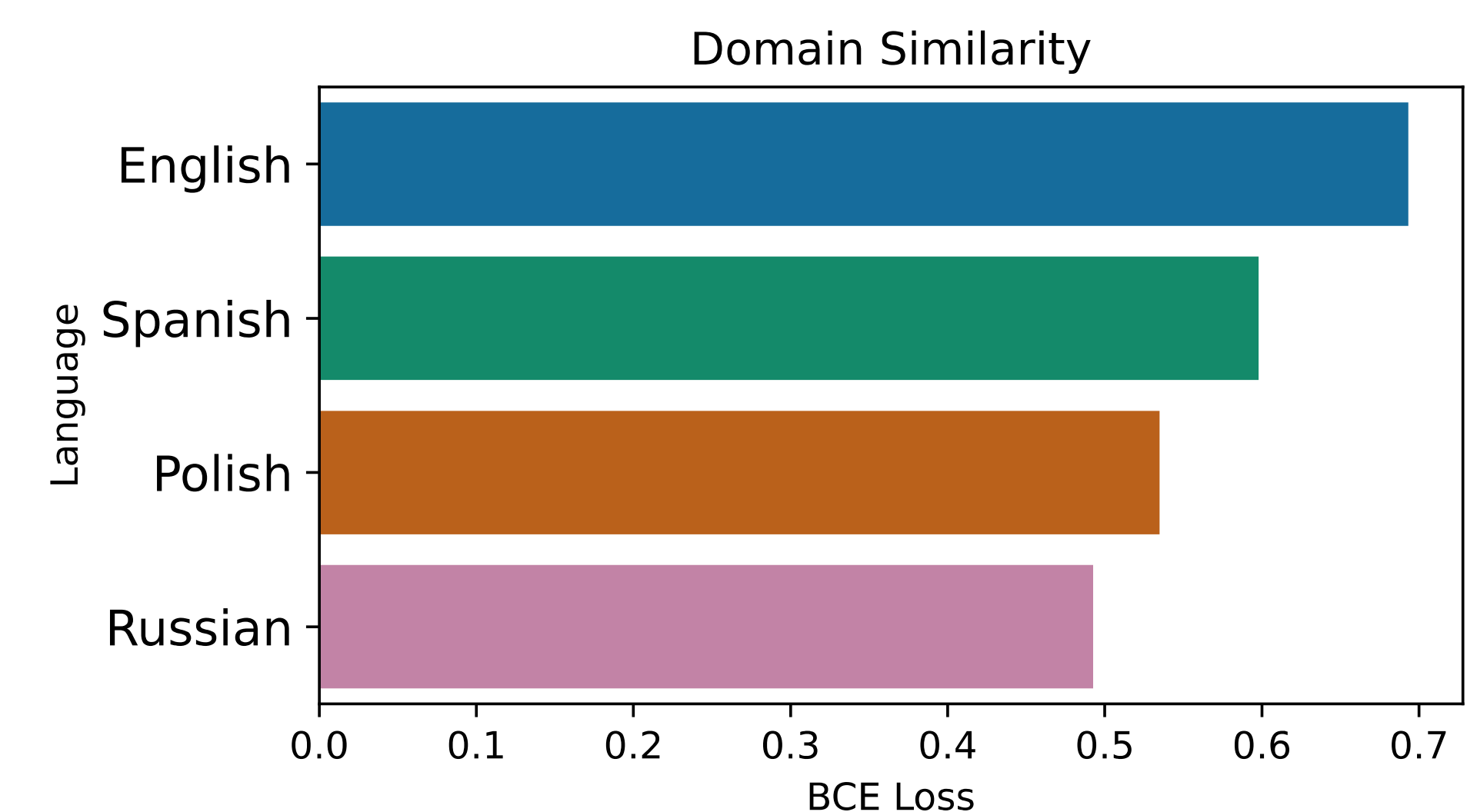
$$ds_{f_1, f_2}^{m_1, m_2} = s(K_{f_1}^{m_1} - \alpha_1 K^0, K_{f_2}^{m_2} - \alpha_2 K^0),$$

where α_i minimizes $\|K_{f_i}^{m_i} - \alpha_i K^0\|_F$.

EXPERIMENTS:

Transfer Learning:

Domain similarity is inconsistent with CKA but consistent with the deconfounded CKA.



[1] Kornblith, Simon, et al. "Similarity of neural network representations revisited." *ICML*, 2019.