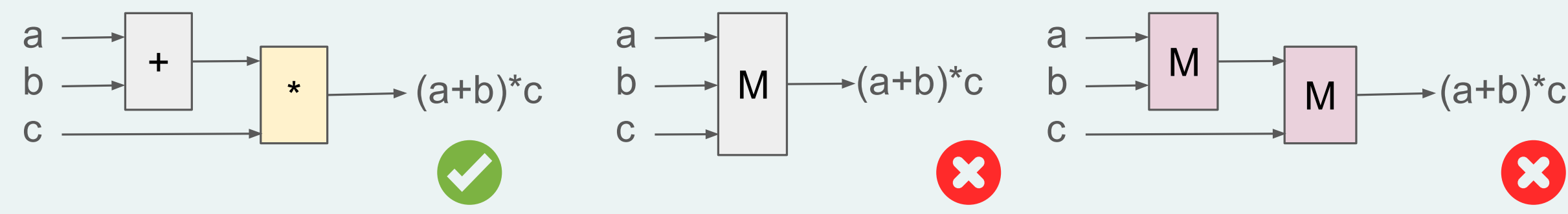


Are Neural Nets Modular? Inspecting Functional Modularity Through Differentiable Weight Masks

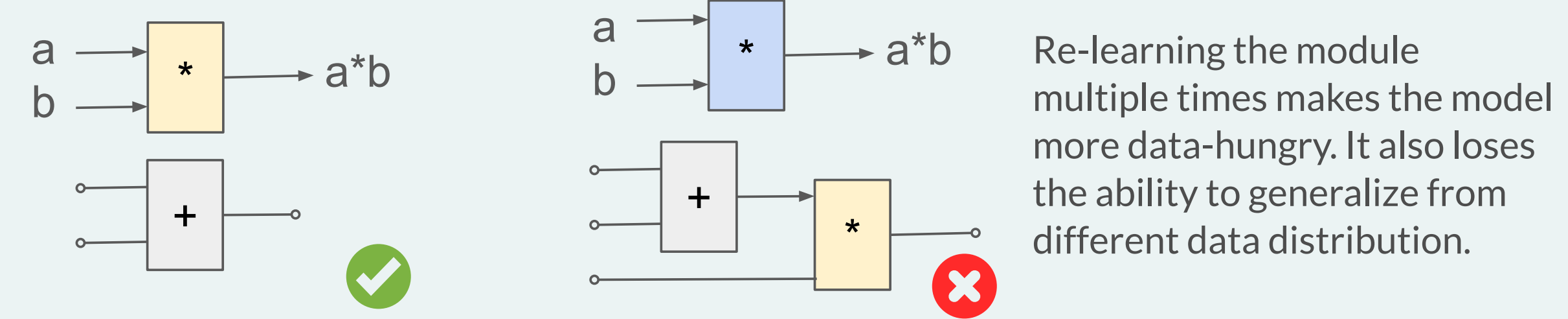
Motivation

- Modularity provides a natural way of achieving compositionality, which appears essential for systematic generalization.
- We investigate if FNNs, LSTMs, Transformers and CNNs are modular, focusing on two desirable properties:

1 Use different modules for separate functions ($P_{\text{specialize}}$)



2 Use the same module for identical functions (P_{reuse})



Method

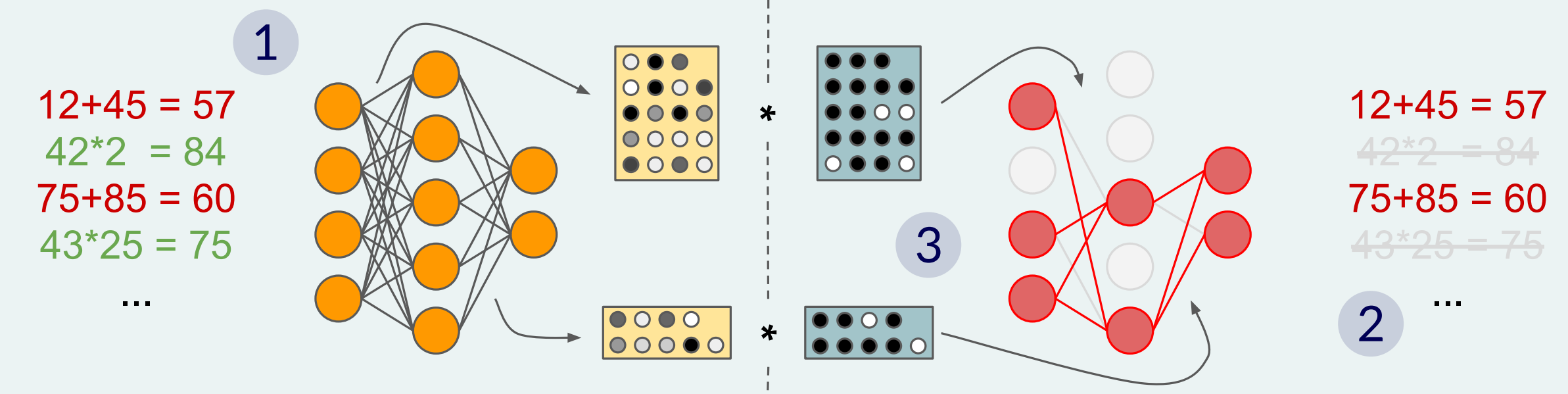
We identify weights responsible for specific functions as follows:

- Pre-train the network on the full task and freeze its weights.
- Create a new dataset focusing on the target functionality. It is usually a subset of the original task.
- Train regularized binary mask for each weight. The weights not needed for the task will be masked out.

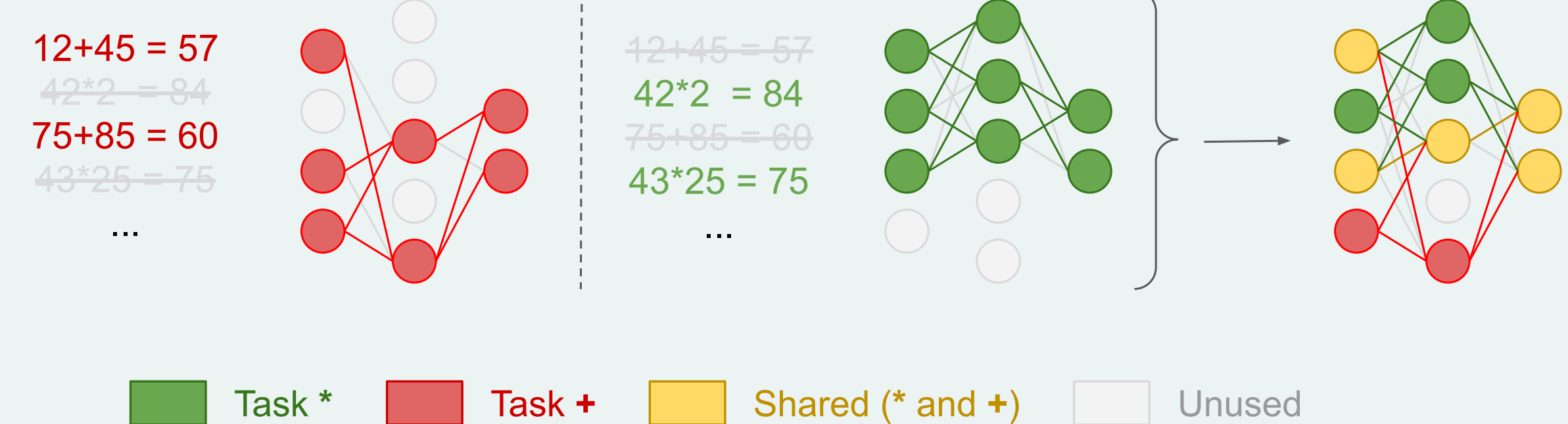
The differentiable binary masks are learned using Gumbel-Sigmoid:

$$s_i = \sigma((l_i - \log(\log U_1 / \log U_2)) / \tau) \quad \text{with } U_1, U_2 \sim U(0, 1)$$

$$w'_i = w_i * b_i, \quad b_i = [\mathbb{1}_{s_i > 0.5} - s_i]_{\text{stop}} + s_i$$



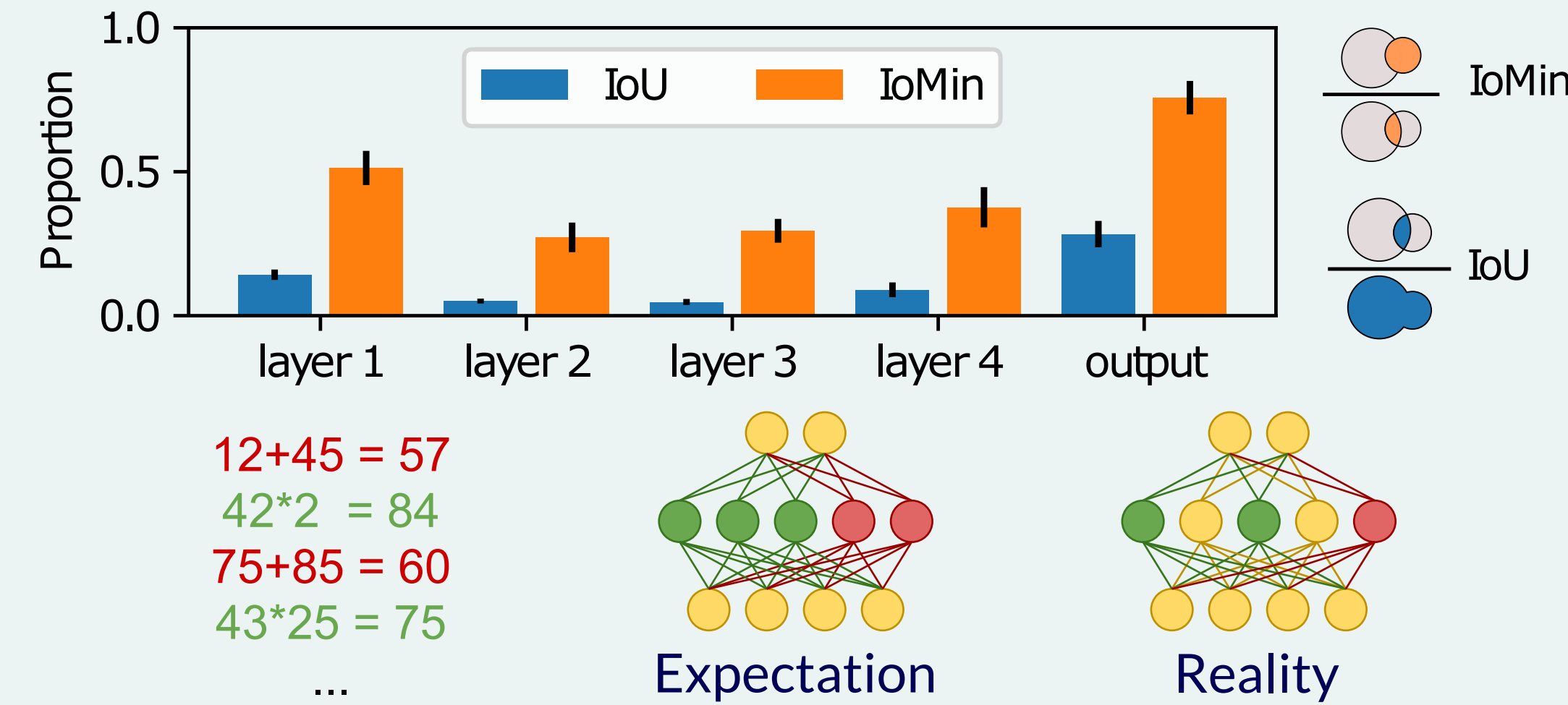
- Repeat the process for all functionalities of interest. This assigns each weight to some functionality.



NNs specialize, but don't reuse

$P_{\text{specialize}}$: Use different modules for separate functions

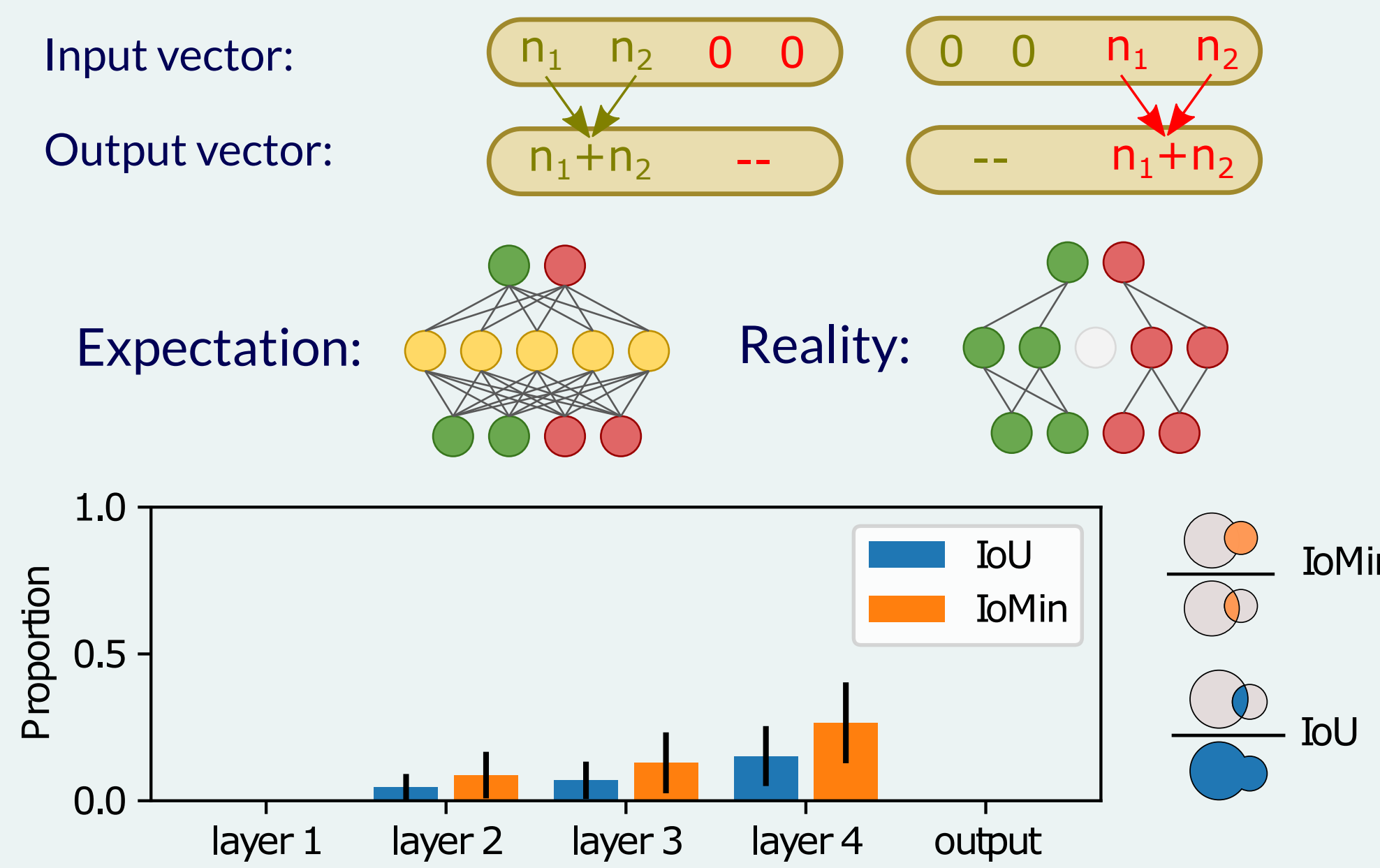
Dataset consisting of two very different functions: addition and multiplication. The inputs are shared between the functions.



There is some separation, but it is not perfect

P_{reuse} : Use the same module for identical functions

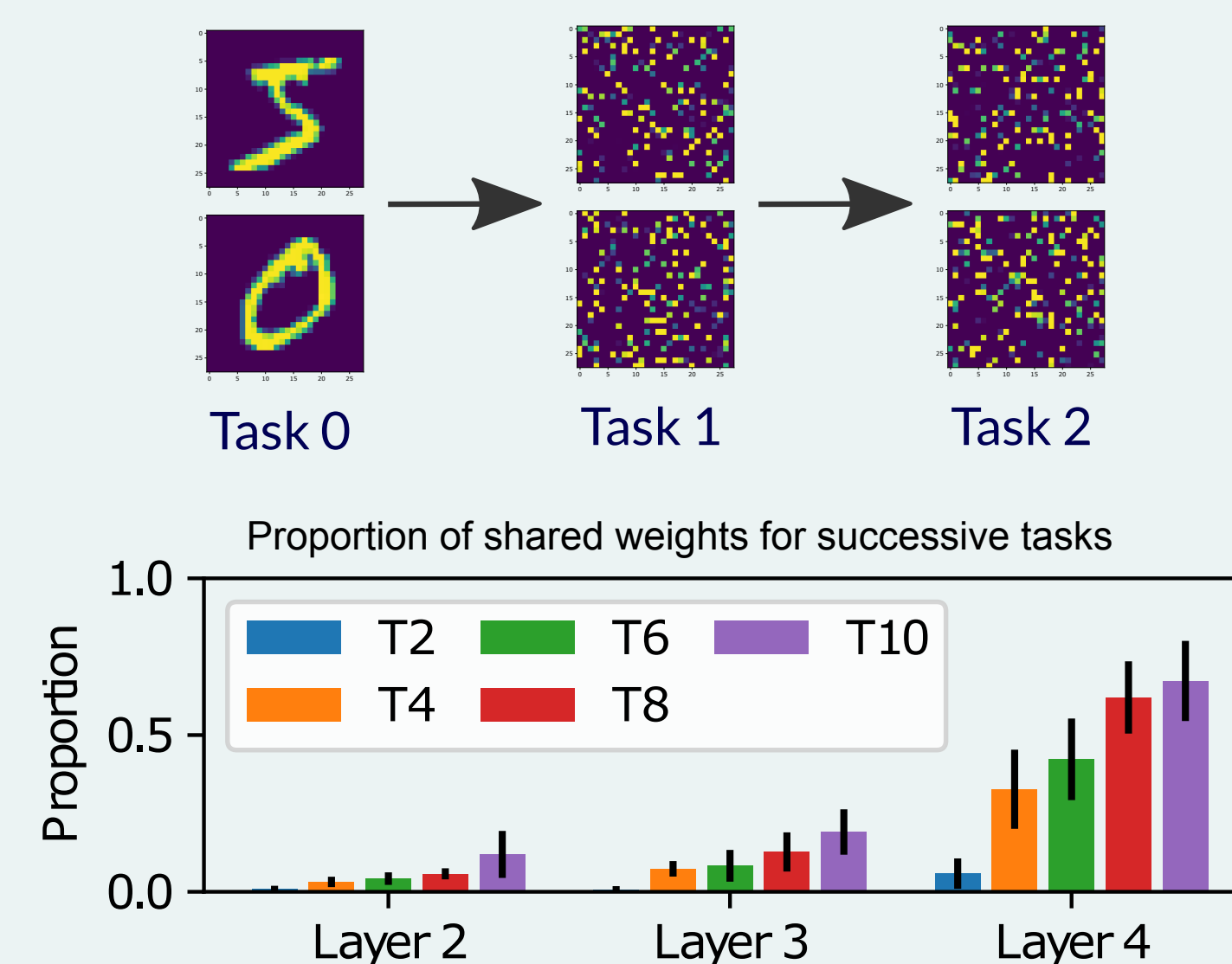
Dataset consisting of the same operation twice. The inputs are not shared.



There is no reuse

Case studies

Continual learning on Permuted MNIST



No weight sharing up until there are enough free weights

- FNN learns the permutations sequentially.
- Weights frozen before learning a new permutation.
- First layer always reinitialized.
- Expectation: re-learn the first layer only, reuse the rest.

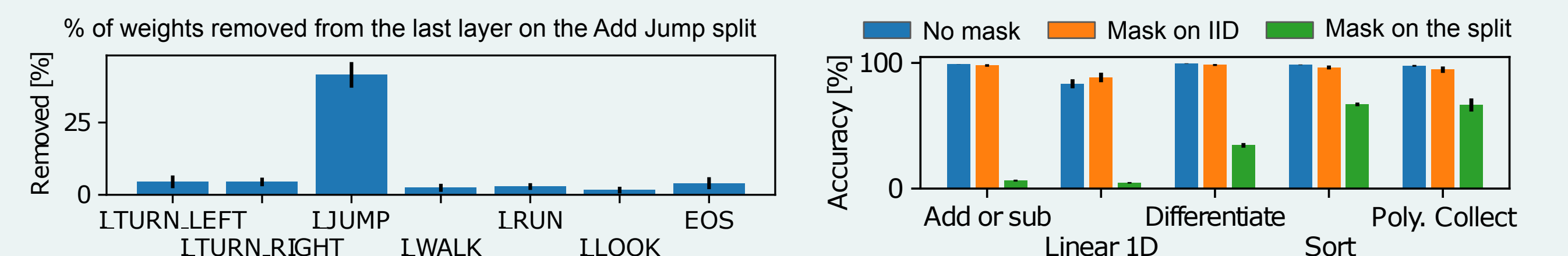
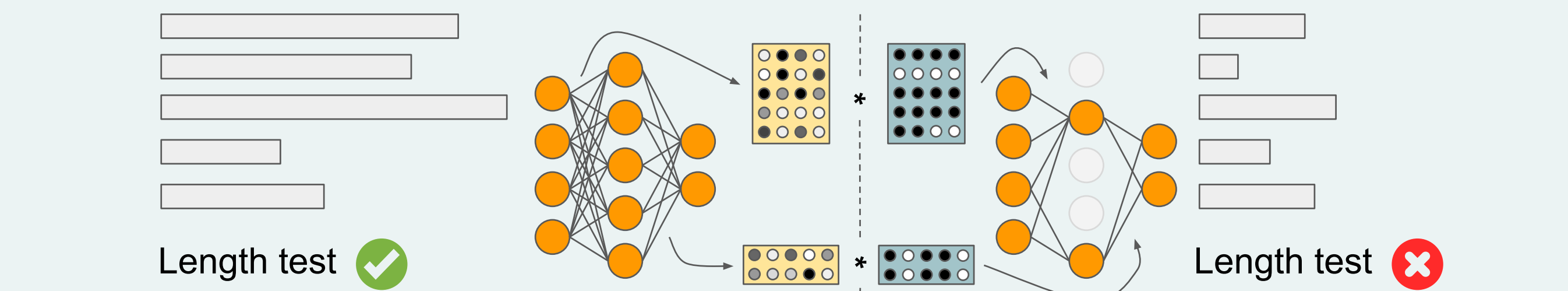
NNs on Algorithmic tasks learn composition specific weights and do not generalize systematically

Which of the following hypothesis explains the generalization issues on SCAN?

- The training set is not enough to find analogies between working and non-working primitives, but the learned algorithm is otherwise general.
- The NN has not learned a general algorithm.

How to find it out?

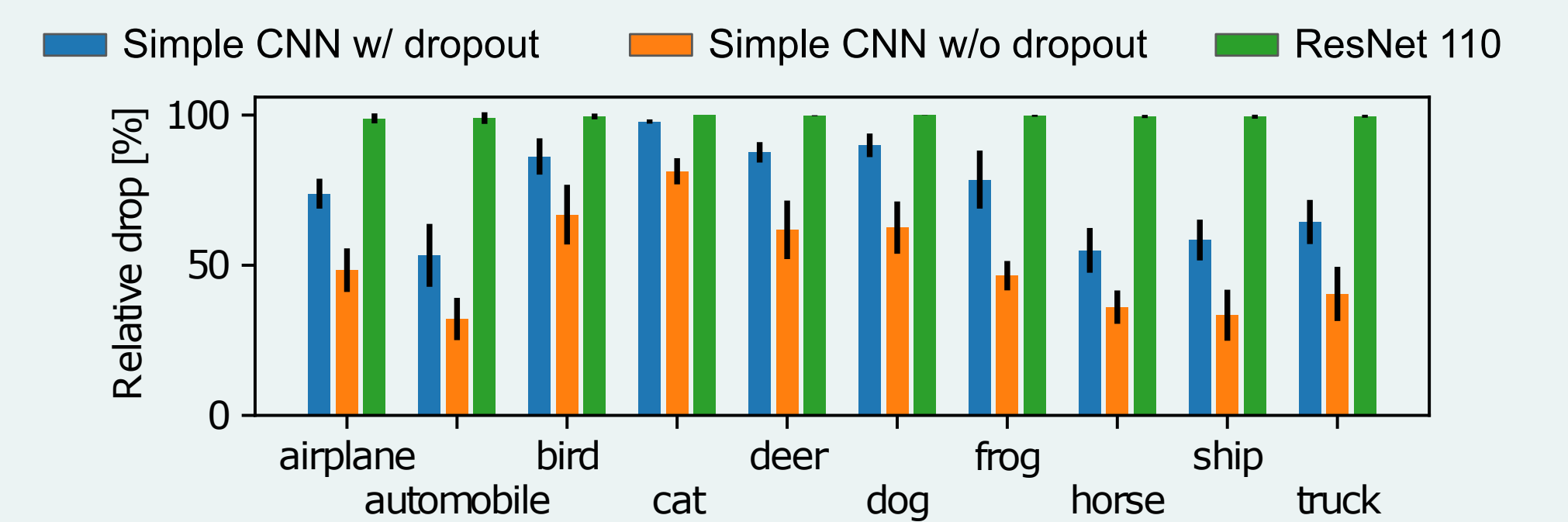
- Pre-train the network on all data. The resulting network performs well on all splits. This rules out **a** as the only issue.
- Train the masks on the specific data splits. We found split-specific weights, meaning the learned algorithm is not general.



This occurs on SCAN & DeepMind Mathematics with both LSTM and Transformers

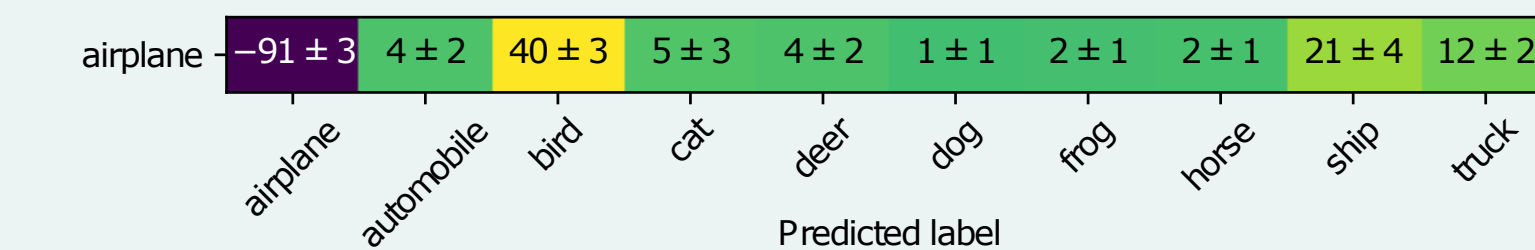
The learned solution is inherently non-compositional

CNNs on CIFAR 10 rely heavily on class-exclusive weights



Removing such weights hurts performance

Analyzing the changes in the confusion matrix reveals class similarities



Conclusion

- ✓ NNs exhibit some level of modularity: different weights become responsible for different functions.
- ✗ NNs resist weight sharing even when the underlying functionality is shared.
- ✗ NNs learn composition-specific weights: thus the learned solution is not general.