# Learning Useful Representations of Recurrent Neural Network Weight Matrices

*Vincent Herrmann, Francesco Faccio, Jürgen Schmidhuber*

IDSIA

Paper

Code & Datasets

Università della Svizzera italiana

Scuola universitaria professionale della Svizzera italiana
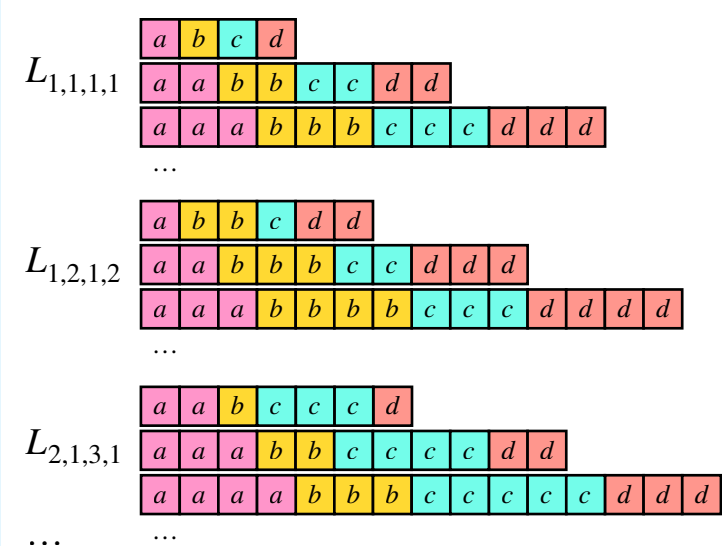
SUPSI

KAUST

## Datasets

- Two datasets of LSTM weights
- Each LSTM is trained to achieve a different task
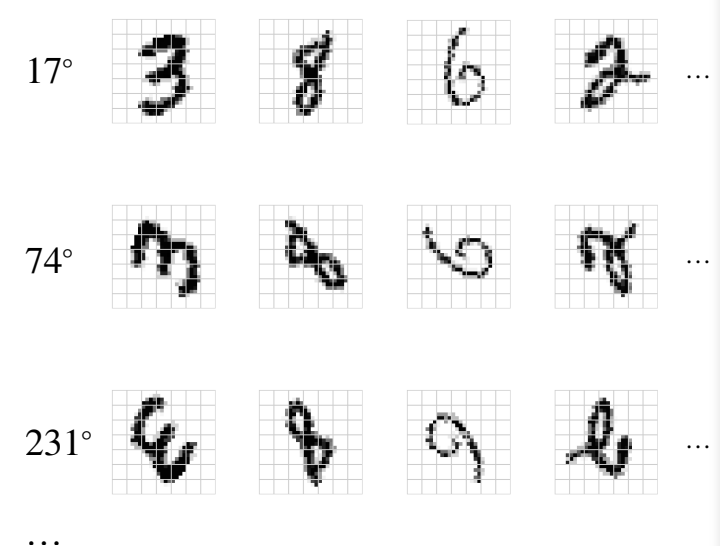
### Formal Languages Dataset

Autoregressive models of languages

$L_{m_a,m_b,m_c,\ldots} := \{a^{n+m_a}b^{n+m_b}c^{n+m_c}\ldots \mid n \in \mathbb{N}\}$

$L_{1,1,1,1}$

$L_{1,2,1,2}$

$L_{2,1,3,1}$
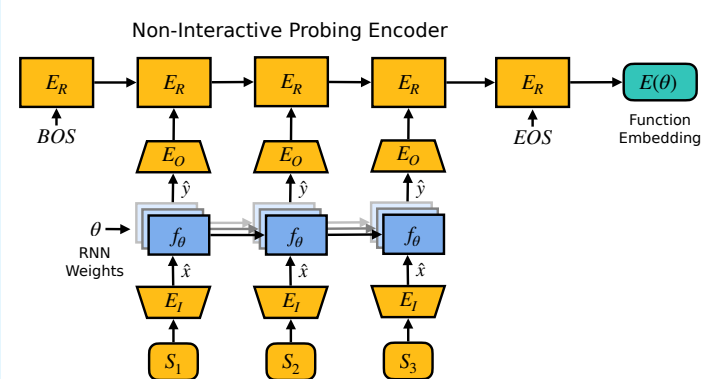
### Tiled Sequential MNIST Dataset

Classifiers of the MNIST dataset, rotated by different angles

17°  74°  231°

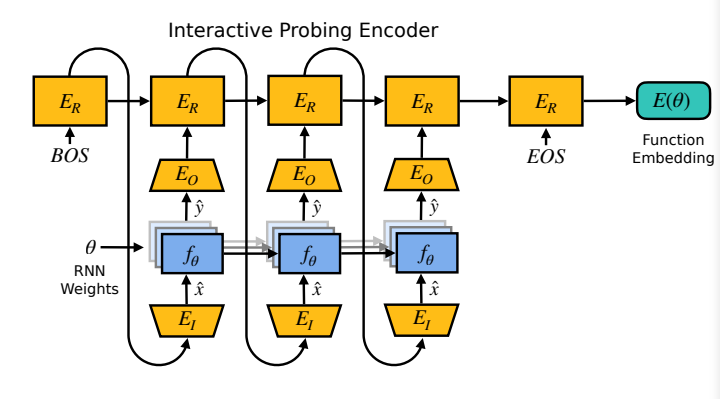## Functionalist Approaches

### Non-Interactive Probing Encoder

- Fixed but learnable probing sequences are given as input to the input RNN $f_\theta$
- Based on the corresponding output sequences, the core LSTM $E_R$ computes the representation $E(\theta)$

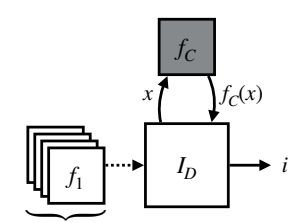Non-Interactive Probing Encoder

### Interactive Probing Encoder

- Probing sequences are dynamically generated by the core LSTM $E_R$
- The next probing input depends on all the previous probing inputs and corresponding outputs

Interactive Probing Encoder

## Theory for the Functionalist Approach

### Setting:

- Interrogator $I_D$ has to identify a specific function $f_C$ from a known set $D$ of total computable functions
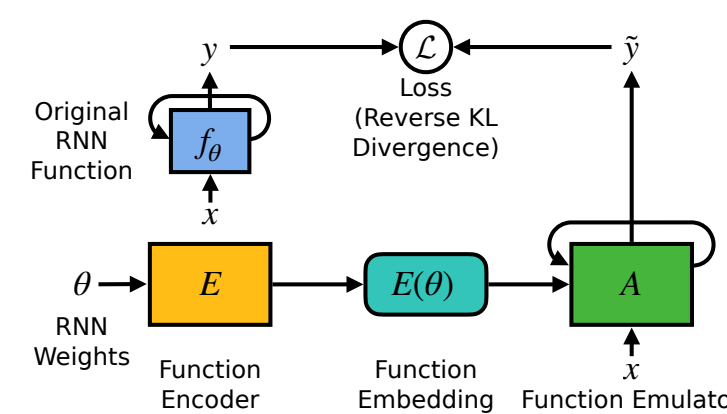- It has to use as few interactions as possible

### Results:

- The general upper bound of required interactions is the same for interactive and non-interactive Interrogators
- For certain function sets, an interactive Interrogator needs exponentially fewer interactions
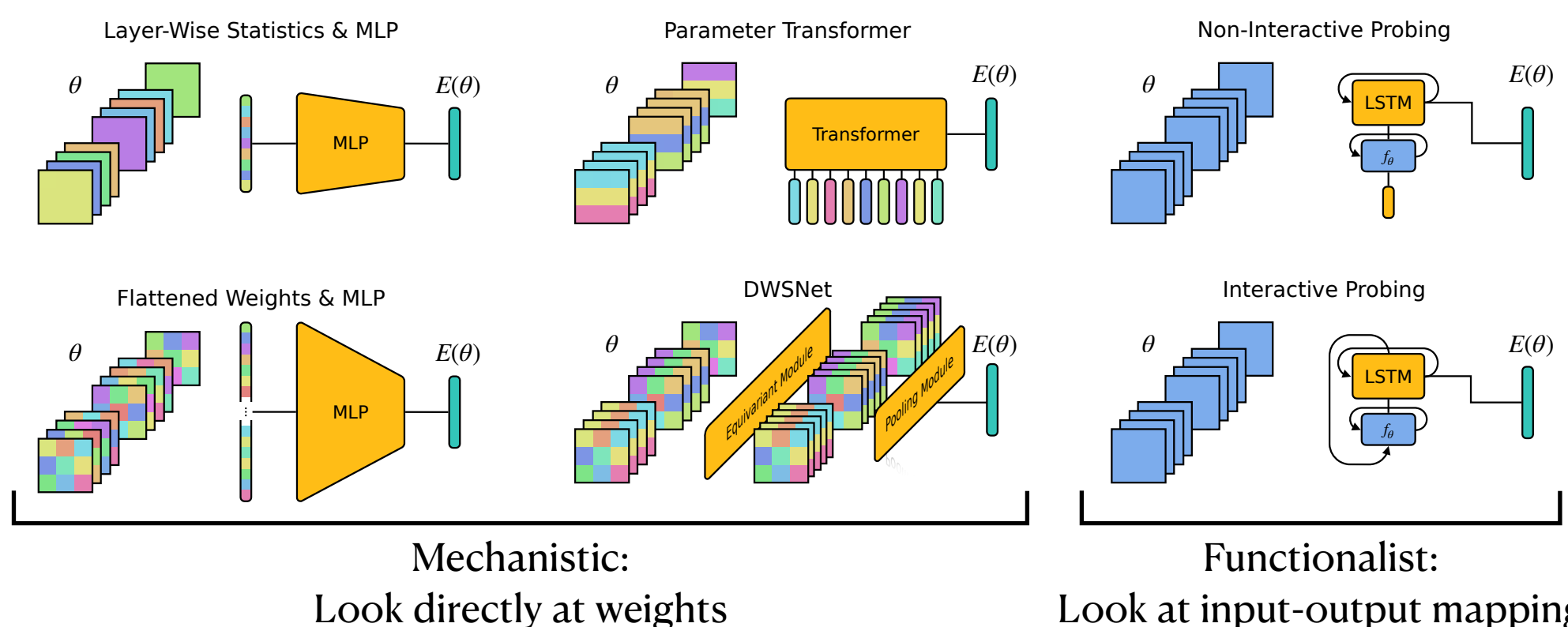
## Recurrent Neural Networks are universal computers. Their weights can represent any program. Can we learn useful representations of the weights of RNNs?

## Self-Supervised Learning of RNN Weight Representations

- Recurrent function $f_\theta$ with parameters $\theta$ is run in an environment, we get a trajectory $S_\theta = (x_1, y_1, x_2, y_2, \ldots)$
- Encoder $E$ generates representation $E(\theta)$
- Emulator $A$ is conditioned on $E(\theta)$ and imitates $f_\theta$

Original RNN Function

Loss (Reverse KL Divergence)

RNN Weights

Function Encoder

Function Embedding

Function Emulator

## Types of Encoders for RNN Weights $\theta$

Layer-Wise Statistics & MLP

$\theta$ → MLP → $E(\theta)$

Parameter Transformer

$\theta$ → Transformer → $E(\theta)$

Non-Interactive Probing

$\theta$ → LSTM / $f_\theta$ → $E(\theta)$

Flattened Weights & MLP

$\theta$ → MLP → $E(\theta)$

DWSNet

$\theta$ → Equivariant Module / Pooling Module → $E(\theta)$

Interactive Probing

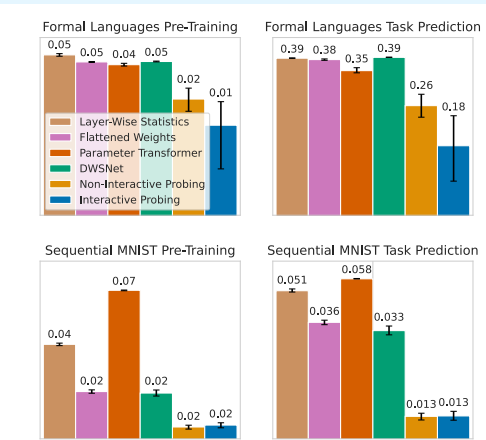$\theta$ → LSTM / $f_\theta$ → $E(\theta)$
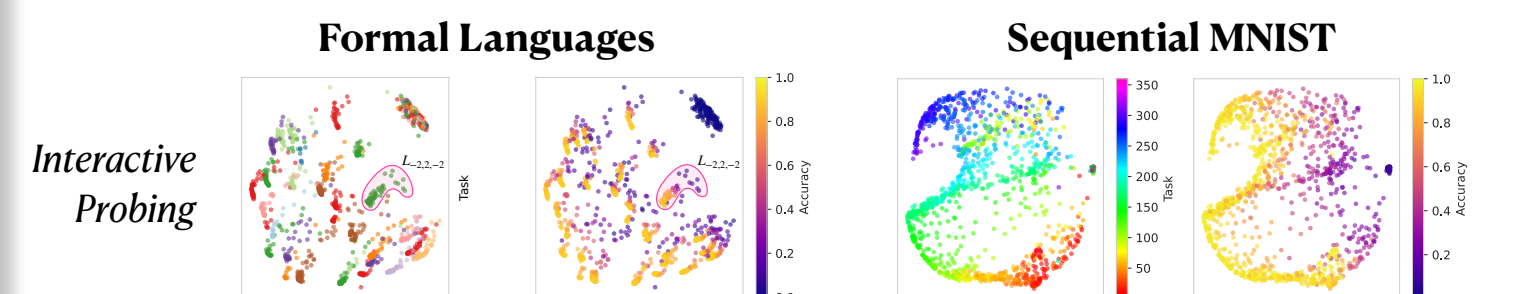
**Mechanistic:** Look directly at weights

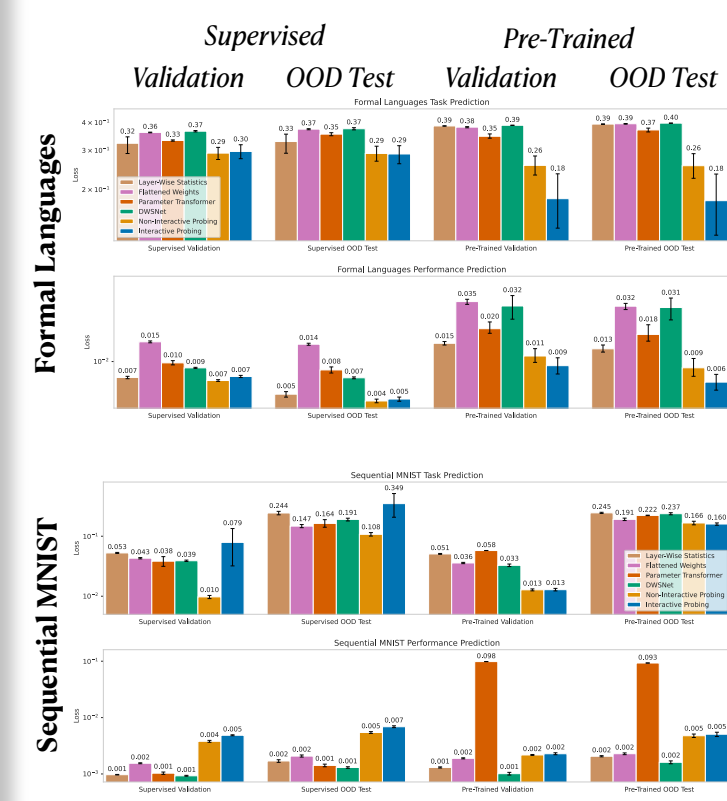**Functionalist:** Look at input-output mapping

## Results

- The learned representations can be used for various downstream tasks, such as task, performance or generalisation gap prediction
- Functionalist approaches are superior at more complex problems
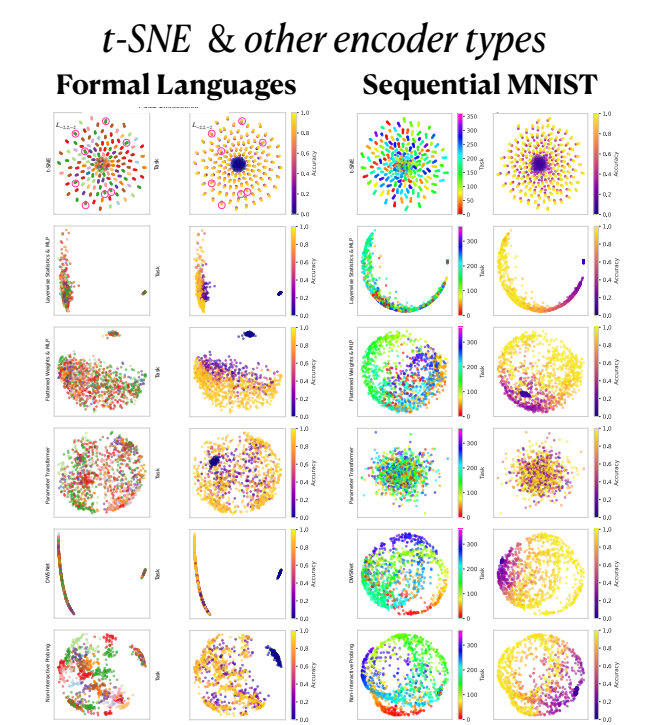- Only Interactive Probing learns generally useful representations for the Formal Languages dataset

Formal Languages Pre-Training

Formal Languages Task Prediction

Sequential MNIST Pre-Training

Sequential MNIST Task Prediction

## Learned Embedding Spaces

Formal Languages

Sequential MNIST

*Interactive Probing*

PCA of learned embedding spaces. Every dot represents a network $f_\theta$ from the validation datasets

## Downstream Results

Supervised Validation | Pre-Trained | Validation | OOD Test

*t-SNE & other encoder types*

Formal Languages

Sequential MNIST

## Original vs. Emulated Performance

$f_\theta$'s original performance vs. the performance of $A_\xi$'s emulation based on $E_\phi(\theta)$. Validation set.

Formal Languages

Sequential MNIST

## Encoder Properties

Invariance to Permutation

*Interactive Probing loss vs. number of probing sequences*

*vs. probing sequence length*

Probing Sequences

Formal Languages Probing Sequences