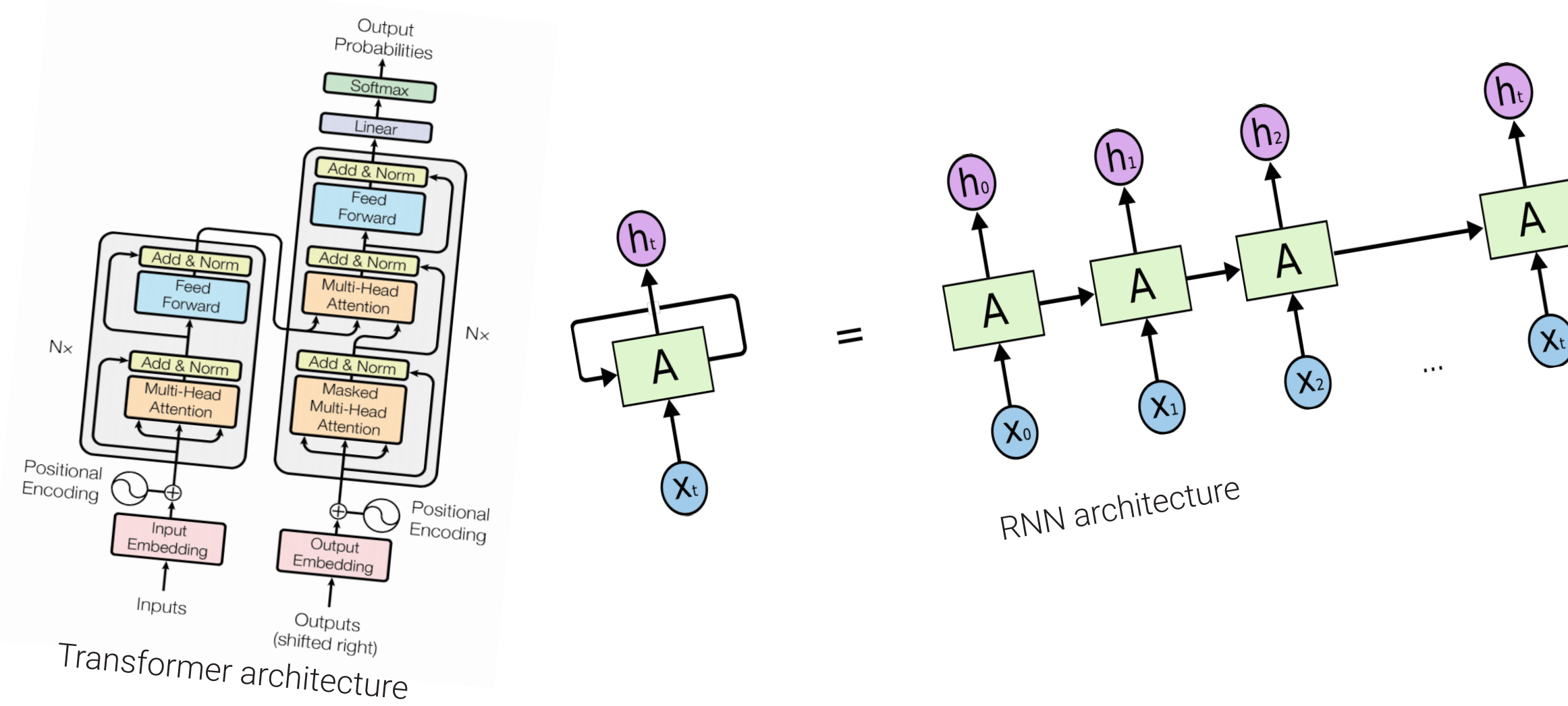


# Neural Network Arena: Investigating Long-Term Dependencies in Deep Models

Hannes Brantner  
Computer Engineering

**TU Wien Informatics**  
Institute of Computer Engineering  
Cyber-Physical Systems Group  
Supervisor: Univ.Prof. Dipl.-Ing. Dr.rer.nat. Radu Grosu  
Contact: e01614466@student.tuwien.ac.at

GitHub repository



## Implemented Models

- LSTM (Long Short-Term Memory)
- GRU (Gated Recurrent Unit)
- CT-RNN (Continuous-Time RNN)
- CT-GRU (Continuous-Time GRU)
- ODE-LSTM (Ordinary Differential Equation LSTM)
- NCP (Neural Circuit Policies)
- Unitary RNN
- Matrix Exponential Unitary RNN **new**
- Unitary NCP **new**
- Transformer
- Recurrent Network Augmented Transformer **new**
- Recurrent Network Attention Transformer **new**
- Memory Augmented Transformer **new**
- Differentiable Neural Computer
- Memory Cell **new**

## Problem Statement

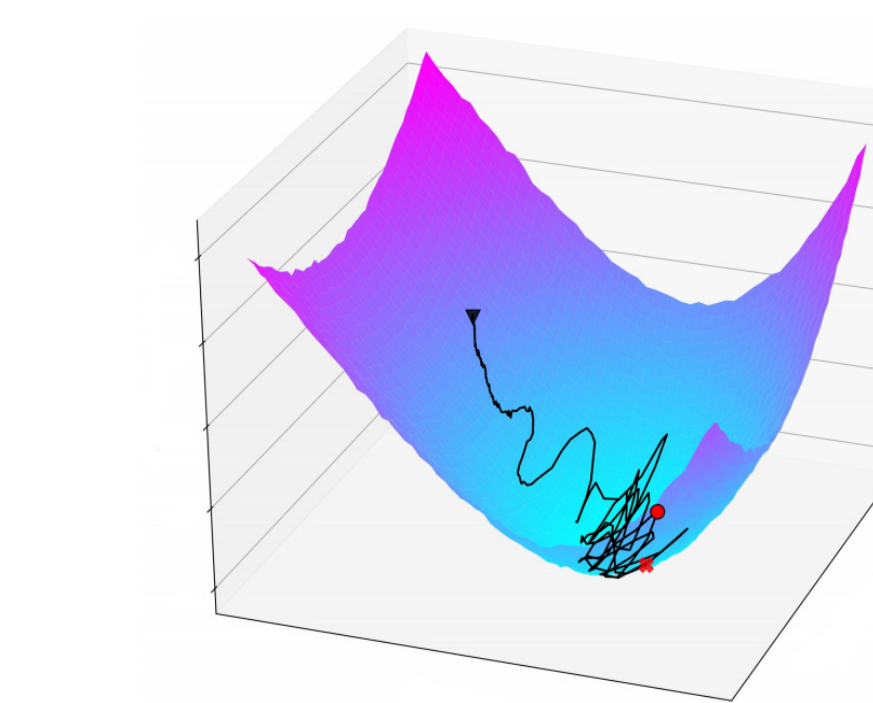
- implementation of a reusable benchmark suite to compare machine learning models used for sequence modeling
- benchmark suite should test the models for their capabilities to capture long-term dependencies and to model physical systems
- selection of state-of-the-art models should be implemented as well as possible outlined improvements
- thoroughful comparison of all implemented models using the benchmark suite
- all implemented models are Transformer or RNN (Recurrent Neural Network) architectures
- proof-of-concept design and implementation of a continuous-time memory cell architecture based on LTC Networks

## Methodology

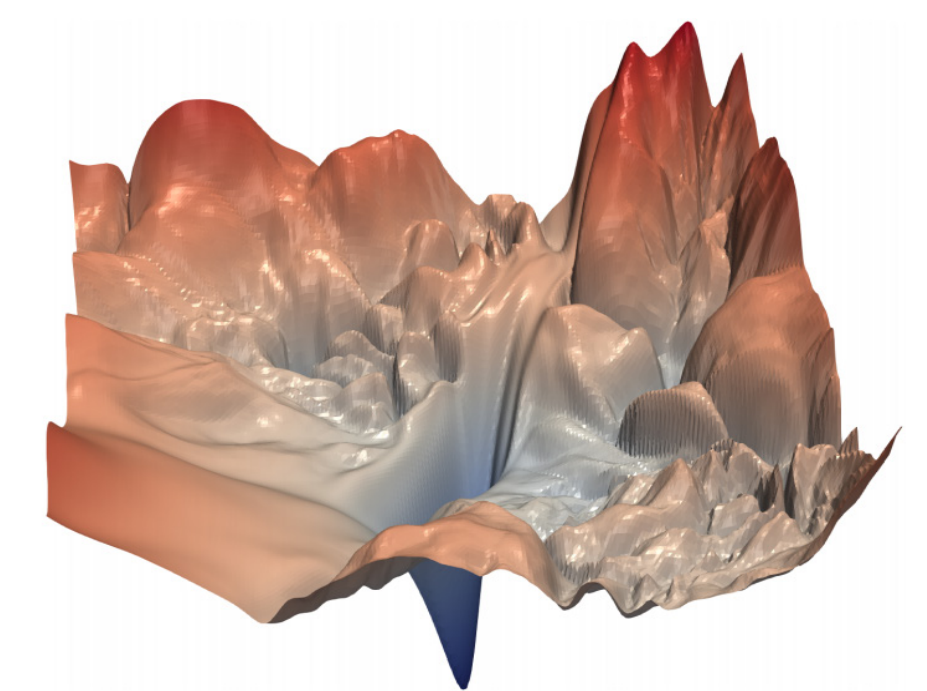
- extensive literature review in the domain of sequence modeling
- implementation of the benchmark suite and all the implemented models
- the benchmark suite was invoked three times on all models and the statistics of the invocation output were interpreted

## Motivation

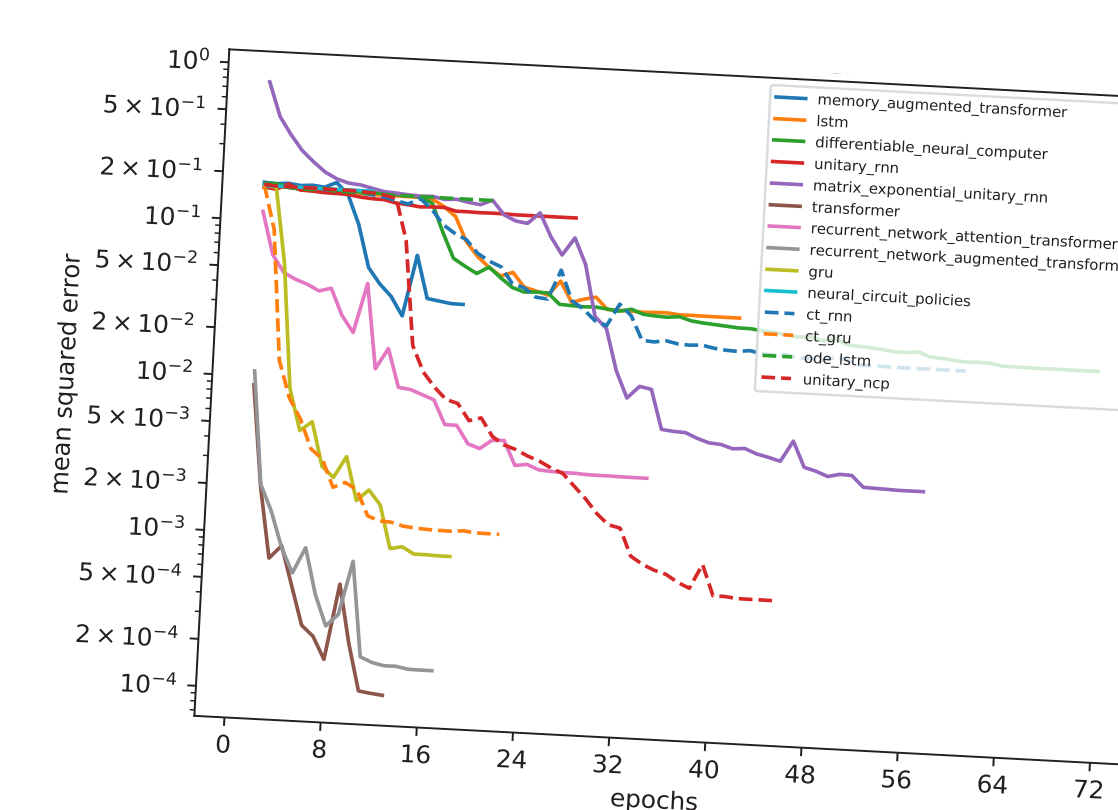
- provide an objective comparison and overview of all implemented models on various sequence modeling tasks
- especially RNN architectures have difficulties of capturing long-term dependencies when being learned by gradient descent
- investigate which mechanism works best in RNN architectures to counteract this difficulty
- each model is just a parameterized function whose parameters are optimized by derivating a loss function that compares model output with expected output
- parameters are updated according to the gradient which should not vanish or explode („walk down the hill“)



loss surface and applied gradient descent procedure



rugged loss surface



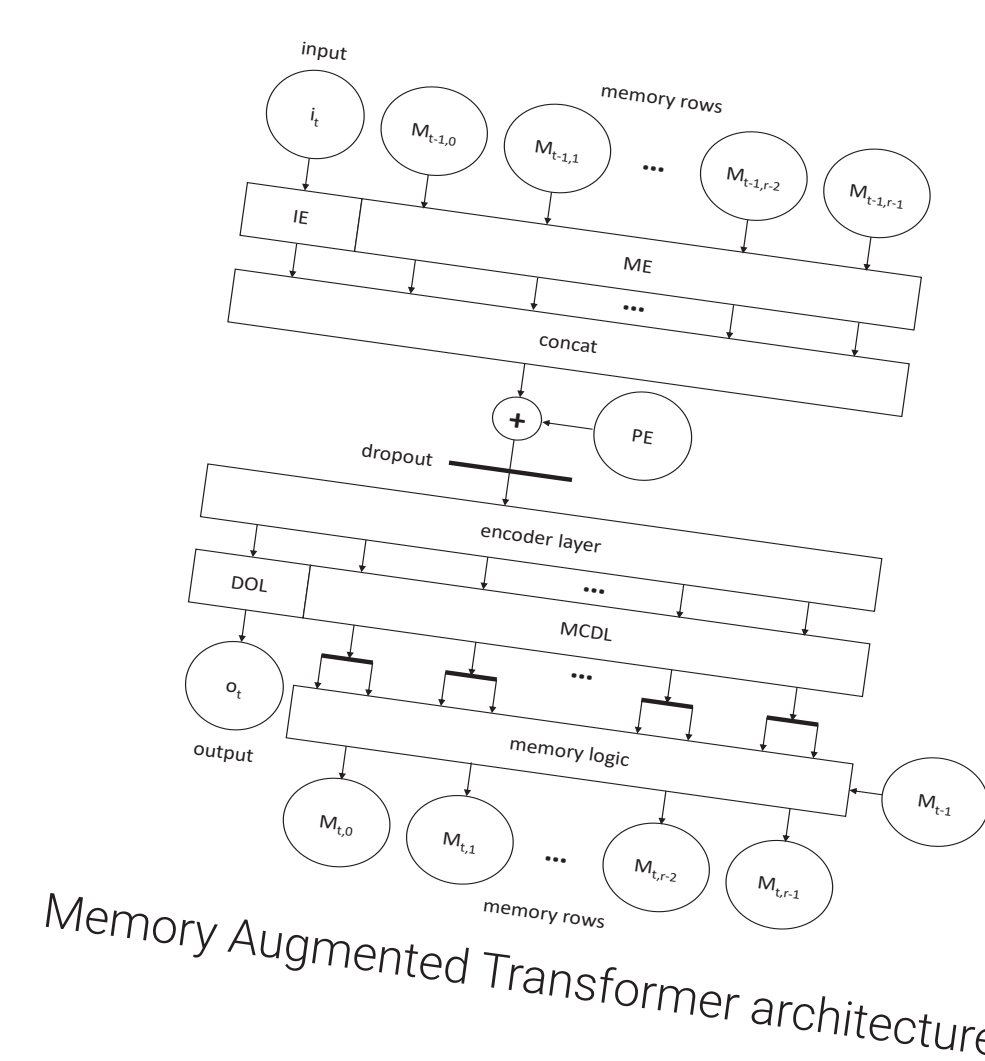
loss evolution during single training run for the Add Benchmark

model	trainable parameters	training duration (s)	training duration per epoch (s)	epochs	test mean squared error
transformer	2189	421.042 ± 62.375	35.175 ± 5.885	12.001 ± 0.237	0.000 ± 0.000
recurrent_network_augmented_transformer	3729	2469.567 ± 238.022	196.200 ± 1.984	17.000 ± 1.000	0.000 ± 0.000
gpr	20041	423.332 ± 32.080	20.821 ± 0.322	20.333 ± 2.517	0.001 ± 0.000
ct_gpr	19073	983.703 ± 123.045	48.464 ± 1.041	20.333 ± 2.517	0.001 ± 0.000
recurrent_neural_network_attention_transformer	2340	1043.323 ± 305.079	39.565 ± 1.507	26.000 ± 7.897	0.001 ± 0.000
recurrent_neural_network_augmented_transformer	21714	613.842 ± 46.210	152.861 ± 2.424	40.333 ± 26.838	0.001 ± 0.000
differentiable_neural_computer	17025	3656.016 ± 153.808	115.400 ± 11.578	31.333 ± 0.028	0.001 ± 0.000
lstm	17469	1556.191 ± 122.046	39.420 ± 0.350	30.333 ± 5.042	0.001 ± 0.000
matrix_exponential_unitary_rnn	17317	986.247 ± 181.783	276.442 ± 11.300	35.000 ± 0.889	0.004 ± 0.076
gtn	2929	1604.727 ± 1483.450	122.520 ± 0.025	32.000 ± 14.731	0.308 ± 0.000
unitary_rnn	1885	3382.152 ± 1513.217	281.514 ± 8.116	14.333 ± 4.628	0.122 ± 0.083
unitary_gnn	1886	3882.347 ± 1304.816	107.356 ± 0.714	15.000 ± 2.000	0.166 ± 0.001
memory_augmented_transformer	1340	748.046 ± 109.977	64.046 ± 0.615	17.000 ± 2.646	0.167 ± 0.001
memory_cell	2537				

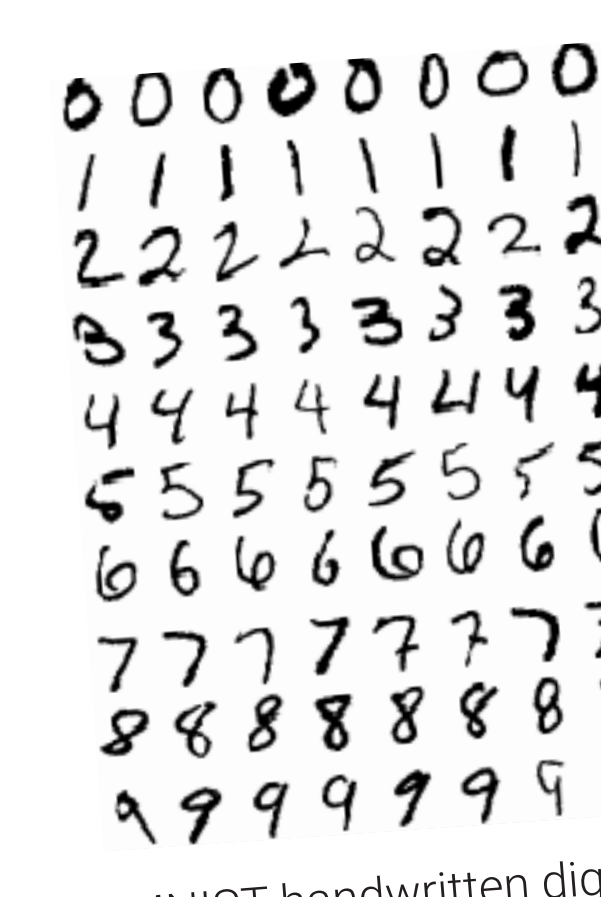
result statistics of Add Benchmark ( $\mu \pm \sigma$ , N = 3)

## Results

- state-of-the-art in efficient unitary matrix parameterization was improved by using an approximated matrix exponential
- continuous-time memory cell architecture was successfully trained to store sparse activations
- positional encoding used in Transformer architectures shows deficiencies in tasks where exact positional information is required
- the newly introduced Memory Augmented Transformer architecture shows promising results in some tasks



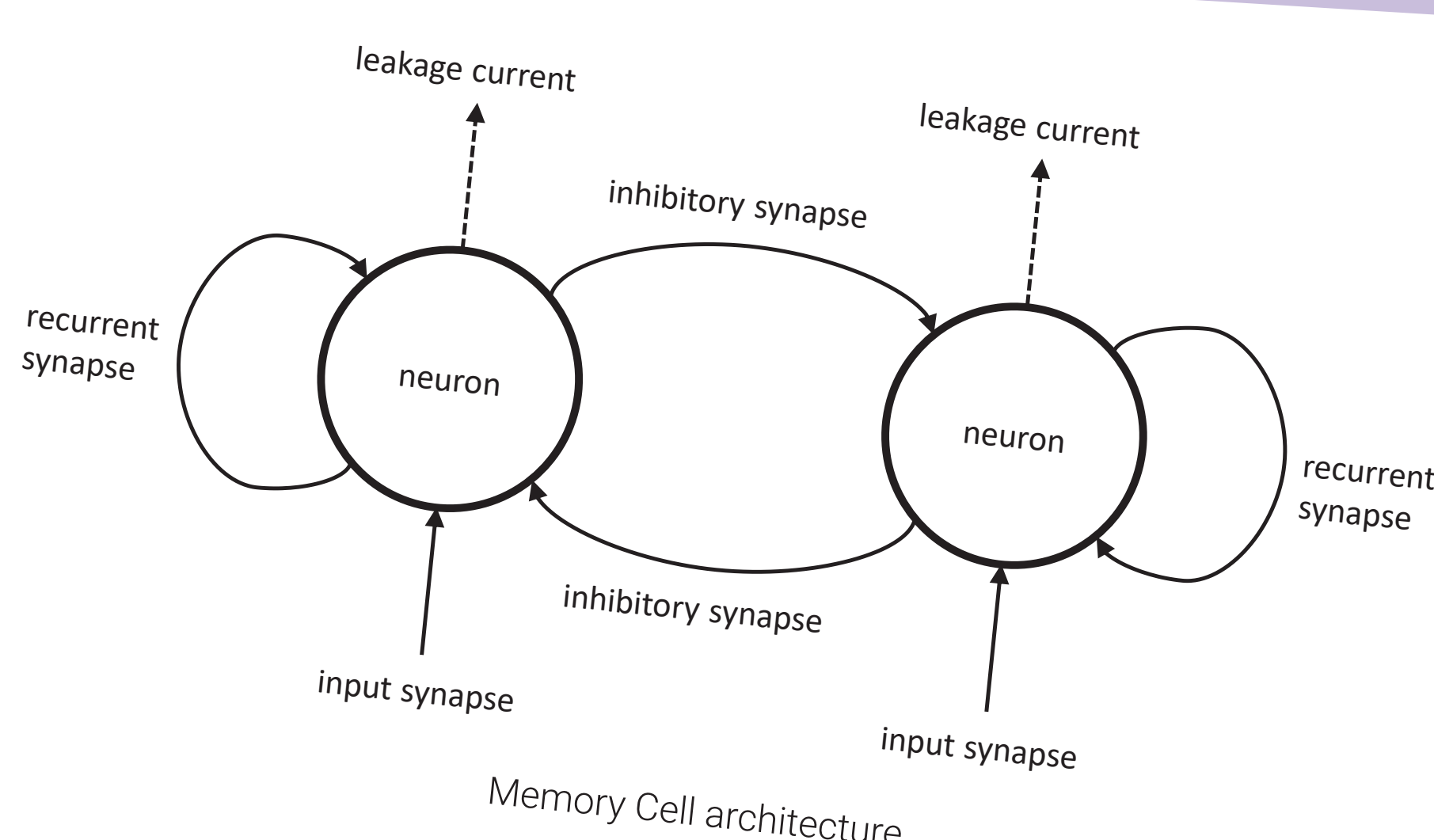
Memory Augmented Transformer architecture



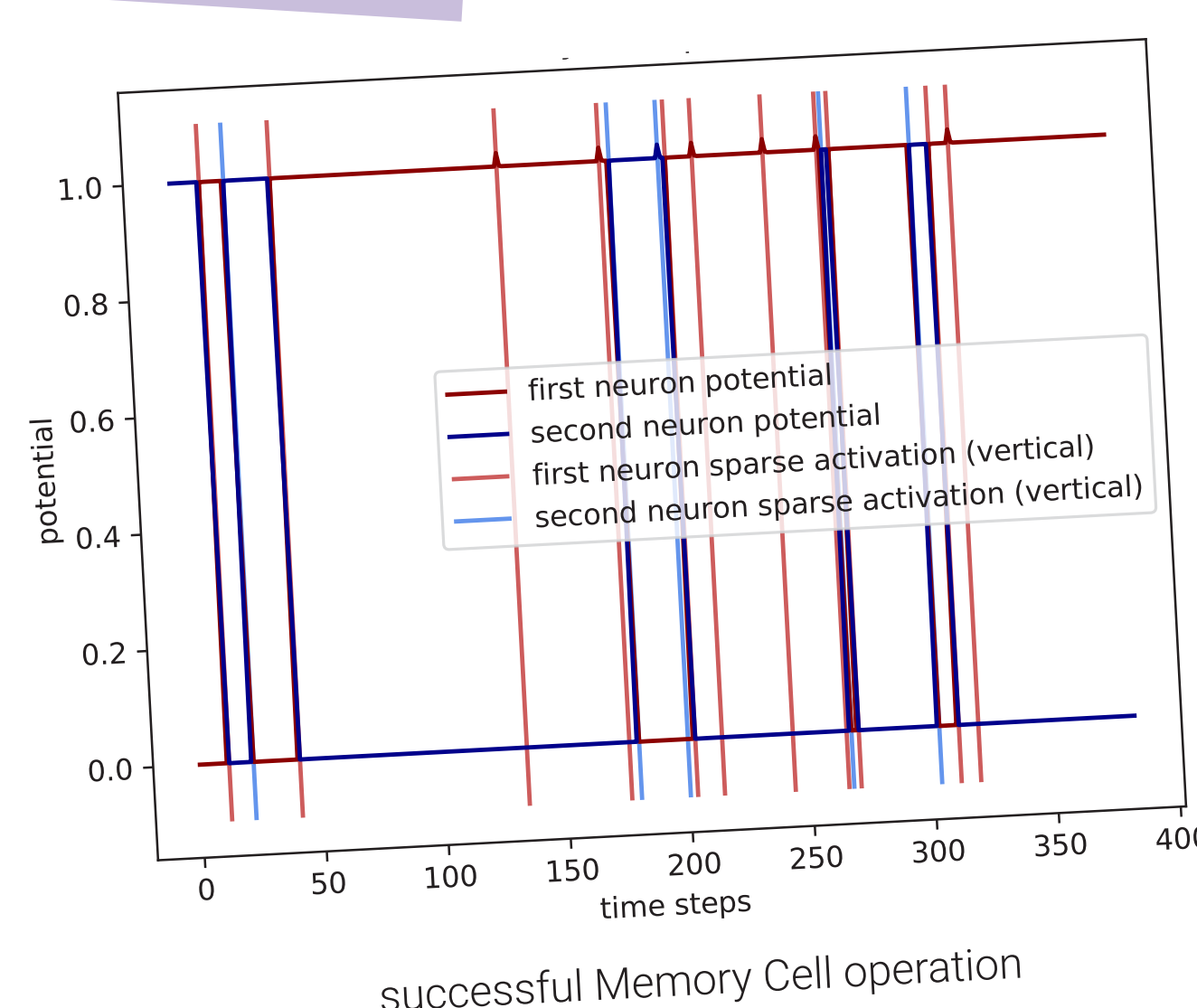
MNIST handwritten digits



physics simulation



Memory Cell architecture



successful Memory Cell operation

## Benchmark Suite Tasks

- Activity Benchmark - human activity classification of inertial sensor measurement data sequences
- Add Benchmark - adding up two marked numbers in a very long number sequence
- Walker Benchmark - predict the next state of a physics simulation given a sequence of previous simulation states
- Memory Benchmark - store a seen category exactly and recall it after seeing a sequence of irrelevant filler symbols
- MNIST Benchmark - digit classification using a sequence of MNIST handwritten digit image chunks
- Cell Benchmark - validates if sparse activations are correctly stored in the time-continuous memory cell architecture