

UKAEA-CCFE-PR(23)187

Wiera Bielajewa, Michelle Tindall, Perumal
Nithiarasu

Comparative study of Transformer- and LSTM-based machine learning methods for transient thermal field reconstruction

Enquiries about copyright and reproduction should in the first instance be addressed to the UKAEA Publications Officer, Culham Science Centre, Building K1/O/83 Abingdon, Oxfordshire, OX14 3DB, UK. The United Kingdom Atomic Energy Authority is the copyright holder.

The contents of this document and all other UKAEA Preprints, Reports and Conference Papers are available to view online free at scientific-publications.ukaea.uk/

Comparative study of Transformer- and LSTM-based machine learning methods for transient thermal field reconstruction

Wiera Bielajewa, Michelle Tindall, Perumal Nithiarasu

Comparative study of Transformer- and LSTM-based machine learning methods for transient thermal field reconstruction

Wiera Bielajewa¹, Michelle Tindall², Perumal Nithiarasu^{1*}

¹Zienkiewicz Institute for Modelling, Data and AI, Swansea University, Bay Campus, Swansea, SA1 8EN UK.

²Culham Science Centre, United Kingdom Atomic Energy Authority (UKAEA), Abingdon, OX14 3DB UK.

*Corresponding author(s). E-mail(s): p.nithiarasu@swansea.ac.uk;
Contributing authors: 946284@swansea.ac.uk; michelle.tindall@ukaea.uk;

Abstract

Solution reconstruction from limited number of measurements is useful in many areas of heat transfer applications. Unlike the standard problems, such reconstruction problems are ill-posed; thus, the non-uniqueness of solution and inherent instability severely complicates the modelling process. Consequently, more conventional inverse analysis methods to reconstruct solutions remain computationally intractable and lacking sufficient flexibility, especially when dealing with time-dependent problems. Aided by powerful Graphical Processing Units (GPUs), Machine Learning (ML) methods rose in popularity due to their flexibility and ability to efficiently process large amounts of data. In recent years, the Transformer-based ML models have gained recognition for their remarkable performance in Natural Language Processing (NLP) tasks as well as time-series analysis, overshadowing the performance of the ML models conventionally used for sequence processing, such as the long short-term memory (LSTM) models. These achievements make Transformer-based models seemingly ideal candidates for reconstructing full solutions from a few measurements. This article compares the performance of these novel Transformer-based models with a simple LSTM model in reconstructing transient one-dimensional (1D) and two-dimensional (2D) thermal fields using sparse spatial measurements. Counterintuitively, the simple LSTM model achieves higher or comparable prediction accuracy compared to the complex Transformer-based models while also exhibiting shorter or comparable training times, which may render Transformer-based models a suboptimal choice for reconstructing transient solutions. Instead, more traditional sequence processing ML models, such as LSTM, might be preferred for this purpose.

Keywords: Machine Learning, Transformer, Transient problem, Solution reconstruction, Conduction, Computational heat transfer, Sparse measurements

Nomenclature

α	Thermal diffusivity
α_x	Thermal diffusivity in x direction
α_y	Thermal diffusivity in y direction
$[\mathbf{K}]$	Key matrix for self-attention operation
$[\mathbf{Q}]$	Query matrix for self-attention operation
$[\mathbf{U}_{pred}]$	Predicted temperature matrix

$[U_{true}]$	True temperature matrix
$[V]$	Value matrix for self-attention operation
$[W]'$	Scaled weight matrix for self-attention operation
$[W]$	Normalised weight matrix for self-attention operation
$[W_k]$	Key weight matrix for self-attention operation
$[W_q]$	Query weight matrix for self-attention operation
$[W_v]$	Value weight matrix for self-attention operation
$[X]$	Input matrix for self-attention operation
$[Y]$	Output matrix for self-attention operation
$\{k_i\}$	Key vector for self-attention operation
$\{q_i\}$	Query vector for self-attention operation
$\{v_i\}$	Value vector for self-attention operation
$\{x_i\}$	Input vector for self-attention operation
$\{y_i\}$	Output vector for self-attention operation
b	Bias of the RNN cell
$b_f, b_i, b_{\bar{c}}, b_0$	Biases of the modified LSTM cell
$b_i, b_{\bar{c}}, b_o$	Biases of the original LSTM cell
c_t	Original LSTM cell state at time t
d_{model}	Model dimension
f_t	Forget gate value of the modified LSTM cell at time t
h	Number of heads in multi-head (self-)attention
h_t	Hidden state (or recurrent information) of the RNN, original and modified LSTM cells at time t
it	Iteration number
l	Prediction window size (sequence length)
lr	Learning rate
N_{out}	Number of output channels
N_{total}	Total number of temperature data points
$PE_{i,n}$	Prediction error at node n at time step i
t	Time
u	Temperature
$u_{pred.i,n}$	Predicted temperature at node n at time step i
$u_{pred.k}$	k^{th} temperature value predicted by the ML model
$u_{true.mean}$	True temperature mean
$u_{true.i,n}$	True temperature at node n at time step i
$u_{true.k}$	k^{th} true temperature value out of N_{total} temperature data points
W_h, W_x	Weights of the RNN cell
$W_{fh}, W_{fx}, W_{ih}, W_{ix}, W_{\bar{c}h}, W_{\bar{c}x}, W_{oh}, W_{ox}$	Weights of the modified LSTM cell
$W_{ih}, W_{ix}, W_{\bar{c}h}, W_{\bar{c}x}, W_{oh}, W_{ox}$	Weights of the original LSTM cell
x	Spatial coordinate
x_t	Input of the RNN, original and modified LSTM cells at time t
y	Spatial coordinate
y_t	Output of the RNN cell at time t

1 Introduction

Transient inverse analysis is a research area in computational engineering which addresses solving various time-dependent inverse problems, including solution reconstruction from limited number of measurements. An inverse problem significantly differs from a standard forward problem. Generally, transient forward problems are well-posed; therefore, given the appropriate initial and boundary conditions the numerical solution can be calculated with a defined accuracy (Tarantola, 2004). Contrariwise, one type of transient inverse problem involves reconstructing the full data inside a problem domain using the available sparse data (observations or measurements). Unlike the standard forward problem, the inverse problem is ill-posed (Tarantola, 2004); consequently, the non-uniqueness of the solution and inherent instability severely complicates the modelling process oftentimes making it computationally intractable, especially for challenging problems.

In order to combat the limitations of the more conventional approaches to reconstructing transient solutions, this paper explores the use of two types of Machine Learning (ML) models to aid the

process of obtaining solutions within an acceptable margin of uncertainty. With the advent of powerful Graphical Processing Units (GPUs) ML became popular in many areas of science and engineering due to its flexibility and the ability to process vast amounts of data within a feasible timescale. In recent years, the Transformer-based ML models have attained recognition by achieving outstanding performance in various Natural Language Processing (NLP) tasks, as evidenced by the well-known ChatGPT chatbot (Qin et al., 2023), as well as numerous time-series analysis problems (Wen et al., 2023; Zhou et al., 2021; Wu et al., 2021; Zhou et al., 2022; Lim and Zohren, 2021). These successes in the area of temporal sequence transformation make them seemingly ideally suited for transient inverse analysis. This paper compares the performance of the novel complex Transformer-based models with the performance of simple long short-term memory (LSTM) model (Hochreiter and Schmidhuber, 1997), which is the ML model type traditionally used for the sequential data processing (Yu et al., 2019), for the task of reconstructing the transient one-dimensional (1D) and two-dimensional (2D) thermal fields.

The aim of the current work is to showcase the suitability of the increasingly popular Transformer-based models for reconstructing transient solutions, a problem encountered in many industrial applications, against more conventional ML models, such as LSTM. The source codes used to produce the results can be found at <https://doi.org/10.5281/zenodo.8208286>.

2 Background

This section provides the background information and add context necessary for understanding the significance of the study presented in this paper.

2.1 Transient inverse problem

The general definition of a transient forward problem can be given as determining the time-dependent effects of the given causes using the applicable physical model of a system. In order to solve the transient forward problem and obtain a full solution inside a domain using standard methods, the system parameters, the boundary conditions, and the initial conditions of the system should be prescribed. Contrariwise, the transient inverse problems can be divided into two types:

1. Determining the system parameters from the observed causes and effects. This is a classic definition of an inverse problem (Tarantola, 2004).
2. Determining the causes from the observed time-dependent effects. In essence, it is a task of using the available sparse data inside a domain (observations or measurements) to reconstruct the full data solutions.

Various methods have been employed over the years to deal with the inverse problems; however, historically, more attention has been given to the inverse problems falling under the first type. The more traditional methods include functional analytic regularisation as well as the statistical regularisation, with the most well-known example being the Bayesian inversion (Tarantola, 2004; Arridge et al., 2019). A search-and-optimisation-based approach is another way to obtain solutions to transient problems. For instance, Bangian-Tabrizi and Jaluria (2018) solved the inverse 2D natural convection problem in steady state using a Particle Swarm Optimisation (PSO) algorithm (Zhang et al., 2015). Whereas Arridge et al. (2019) and Tamaddon-Jahromi et al. (2020) provide a more comprehensive review of the various methods used to solve inverse problems.

This paper focuses on the transient inverse problem of the second type, the reconstruction of transient thermal fields in particular. Perhaps the most common sources of the sparse data observed inside a domain are the data obtained from physical experiments. Fusion energy technology research facilities, designed to test components' suitability for the extreme environment inside a fusion reactor, regularly encounter transient inverse problems stemming from the sparse experimental data. The HIVE (Heating by Induction to Verify Extremes) experimental facility (Hancock et al., 2018) is an illustrative example of that: inverse analysis has to be performed to reconstruct the full temperature field using the temperature measurements recorded by a few thermocouples. Figure 1 shows an example of the HIVE experimental setup.

2.2 Long short-term memory (LSTM)

Recurrent neural networks (RNNs) are extensively used for various sequential data processing tasks (Yu et al., 2019). RNNs typically consist of a number of standard recurrent cells (Figure 2), the

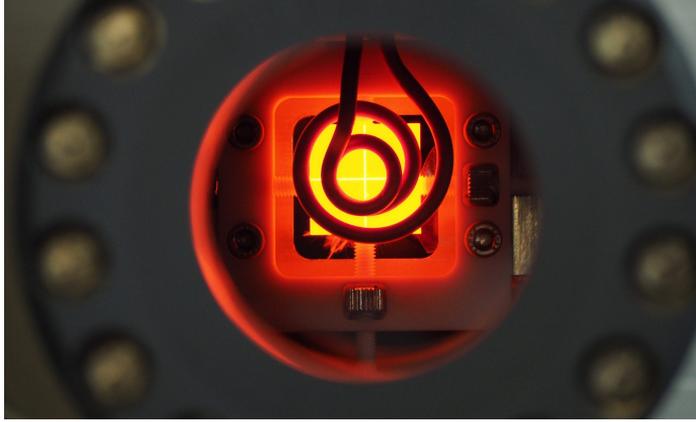


Fig. 1: The arrangement of coil and sample inside HIVE (Hancock et al., 2018).

mathematical representation of which can be written as:

$$\begin{aligned} h_t &= \sigma(W_h h_{t-1} + W_x x_t + b) \\ y_t &= h_t \end{aligned} \quad (1)$$

where x_t , y_t , and h_t represent the input, output, and the hidden state (or recurrent information) of the cell at time t respectively, h_{t-1} is the hidden state at time $t - 1$, while W_h and W_x are the weights of the cell, and b is the bias. The operator σ is a sigmoid function.

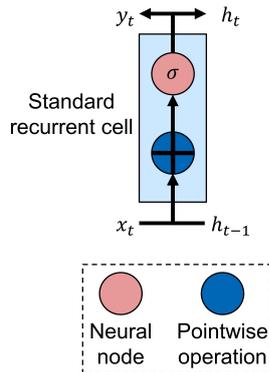


Fig. 2: Standard recurrent cell for RNNs.

However, a RNN comprised of the standard recurrent cells tend to experience difficulties during the training process due to the vanishing or exploding gradients between inputs that are far apart in time (Bengio et al., 1994). In order to address this problem of long-term dependencies, the long short-term memory (LSTM) cell, a type of RNN, was developed more than two decades ago (Hochreiter and Schmidhuber, 1997) and successfully applied to a wide range of sequential tasks, such as speech recognition (Hsu et al., 2016), trajectory prediction (Alché and de La Fortelle, 2017), and prediction of the remaining useful life (Ren et al., 2021), to name a few. The original LSTM cell contains only input and output gates (Figure 3) and is represented by the following expressions:

$$\begin{aligned} i_t &= \sigma(W_{ih}h_{t-1} + W_{ix}x_t + b_i) \\ \tilde{c}_t &= \tanh(W_{\tilde{c}h}h_{t-1} + W_{\tilde{c}x}x_t + b_{\tilde{c}}) \\ c_t &= c_{t-1} + i_t \odot \tilde{c}_t \\ o_t &= \sigma(W_{oh}h_{t-1} + W_{ox}x_t + b_o) \\ h_t &= o_t \odot \tanh(c_t) \end{aligned} \quad (2)$$

where c_t and c_{t-1} represent a LSTM cell state at times t and $t - 1$ respectively, while W_{ih} , W_{ix} , $W_{\tilde{c}h}$, $W_{\tilde{c}x}$, W_{oh} , W_{ox} are the weights, b_i , $b_{\tilde{c}}$, b_o are the biases. The operator \odot is the Hadamard product

(also known as the element-wise product). The input gate determines what information should be stored in the cell state, whereas the output gate chooses what information should be extracted from the cell state for the output.

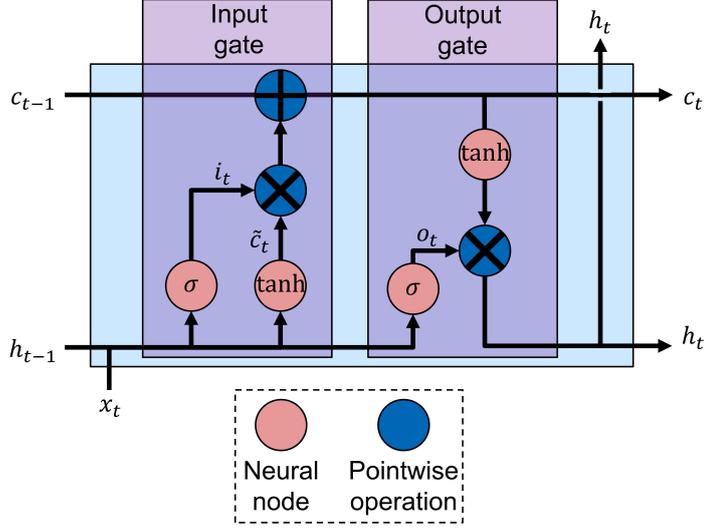


Fig. 3: Original LSTM cell (Hochreiter and Schmidhuber, 1997).

In the present work, a LSTM network comprising a modified version of the original LSTM cells is used. The modification incorporates a forget gate, which decides what information should be eliminated from the cell state (Gers et al., 2000). Figure 4 shows the modified LSTM cell, and the following equations represent the cell:

$$\begin{aligned}
 f_t &= \sigma(W_{fh}h_{t-1} + W_{fx}x_t + b_f) \\
 i_t &= \sigma(W_{ih}h_{t-1} + W_{ix}x_t + b_i) \\
 \tilde{c}_t &= \tanh(W_{\tilde{c}h}h_{t-1} + W_{\tilde{c}x}x_t + b_{\tilde{c}}) \\
 c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \\
 o_t &= \sigma(W_{oh}h_{t-1} + W_{ox}x_t + b_o) \\
 h_t &= o_t \odot \tanh(c_t)
 \end{aligned} \tag{3}$$

where f_t represent the value of the forget gate at a time t .

It should be noted that in this sub-section t is the t^{th} time step in a sequence as it is the commonly used notation for RNNs; however, t means the total sequence length in the subsequent sub-sections.

2.3 Transformers

The classic Transformer model was created by the Google research team in 2017 and successfully applied to the Natural Language Processing (NLP), such as natural language generation and machine translation (Vaswani et al., 2017). Nowadays, Transformer is considered to be the best model for the various NLP tasks (Wolf et al., 2020), with Chat Generative Pre-trained Transformer (ChatGPT) (Qin et al., 2023), an Artificial Intelligence (AI) chatbot, being probably the most famous example of the impressive results Transformers are capable of achieving. However, NLP is essentially a sequence-to-sequence transformation task, since the model input and output are almost always an ordered series of elements. This fact makes the Transformer an appropriate model for time-series prediction. Indeed, in recent years a number of Transformer-based models were successfully applied to various time-series tasks, such as weather, electricity consumption, and exchange rate prediction. Wen et al. (2023) reviewed the state-of-the-art Transformer-based models used to analyse the time series; whereas, Lim and Zohren (2021) conducted a survey of Deep Learning (DL) methods used for time-series forecasting. The aforementioned surveys might not be completely exhaustive; as Transformers are gaining popularity, the new Transformer-based models and the numerous variations of the already existing ones are created every year. Consequently, it is quite challenging to keep track

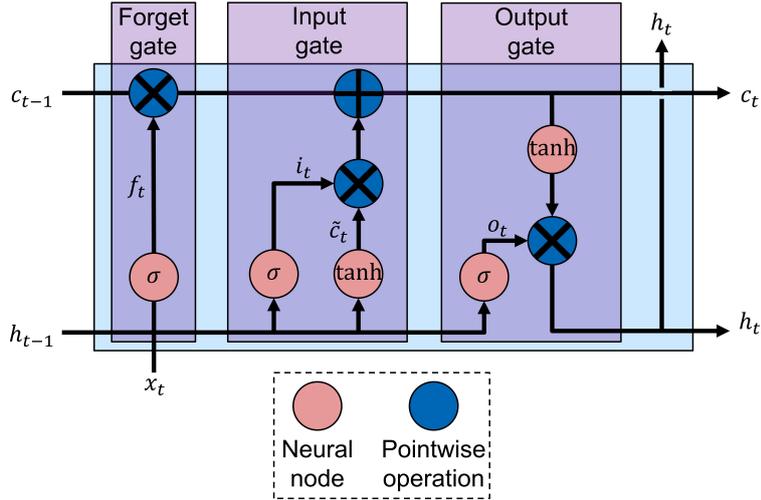


Fig. 4: Modified LSTM cell with a forget gate.

of all new developments, particularly due to the fact that these advances are oftentimes made in widely different areas of research, such as weather prediction (Wu et al., 2022) and computer vision (Ivanovic and Pavone, 2019). Meaning that the process of comparing the new models between each other is complicated by the fact that they have been tested on the research-area-specific datasets, therefore it is not immediately obvious which model is better suited for engineering applications.

The primary advantage of the Transformer-based models lies in the absence of any recurrent connections, which are present in RNNs (Hochreiter and Schmidhuber, 1997). The elimination of the recurrent connections should significantly reduce the training times and make the model more parallelisable, which is beneficial for the model training on GPUs. Furthermore, they were found to be excellent at detecting the long- and short-term sequence dependencies, which should make the model more accurate. (Vaswani et al., 2017).

Figure 5 shows a simplified transformer block. This paper will not provide a detailed explanation of how Transformers work. Readers are referred to Bloem (2019) and Vaswani et al. (2017) for a more detailed introduction to self-attention and the classic Transformer. Nevertheless, a brief overview of the classic self-attention is provided in the next sub-section.

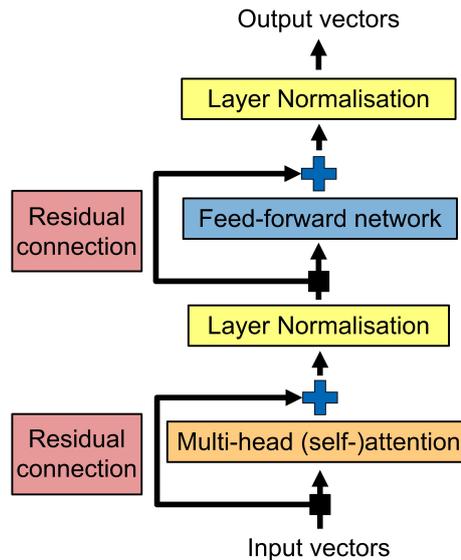


Fig. 5: Simplified transformer block, which consists of multi-head (self-)attention, layer normalisation, feed-forward network, and another layer normalisation. The reader is referred to Bloem (2019) for a more detailed introduction to the classic Transformer.

2.4 Self-attention

At the core of any Transformer-based model is a self-attention operation or a variation thereof. Figure 6 shows the classic self-attention operation used in the classic Transformer (Vaswani et al., 2017). The input is t vectors of size K , while the generated output is a different set of t vectors of the same size. The input vectors are used to calculate queries, keys, and values using query, key, and value weight matrices:

$$\begin{aligned} \{q_i\} &= [W_q] \{x_i\} \\ \{k_i\} &= [W_k] \{x_i\} \\ \{v_i\} &= [W_v] \{x_i\} \end{aligned} \quad (4)$$

where $\{q_i\}$, $\{k_i\}$, $\{v_i\}$, and $\{x_i\}$ are i^{th} query, key, value, and input vectors, respectively; $[W_q]$, $[W_k]$, and $[W_v]$ are the query, key, and value weight matrices, respectively. All query, key, and value vectors can be concatenated to obtain query, key, and value matrices, which are $[Q]$, $[K]$, and $[V]$, respectively. The scaled weight matrix $[W]'$ is calculated using the following equation:

$$[W]' = \frac{[Q][K]^T}{\sqrt{K}} \quad (5)$$

The scaled weight matrix is normalised using a softmax function (Goodfellow et al., 2016):

$$[W] = \text{softmax}([W]') \quad (6)$$

Finally, the normalised weight matrix $[W]$ is then multiplied by the values vector to finally obtain the output matrix $[Y]$, consisting of output vectors $\{y_i\}$:

$$[Y] = [W][V] \quad (7)$$

The query, key, and value weight matrices are trainable parameters. The Transformer-based models tend to employ several self-attention operations in parallel (Figure 7), which allows for a more efficient extraction of the various features in the given time series. The self-attention operation comprises a number of matrix multiplications, which is beneficial as they can be performed using a highly optimised and efficient matrix multiplication code.

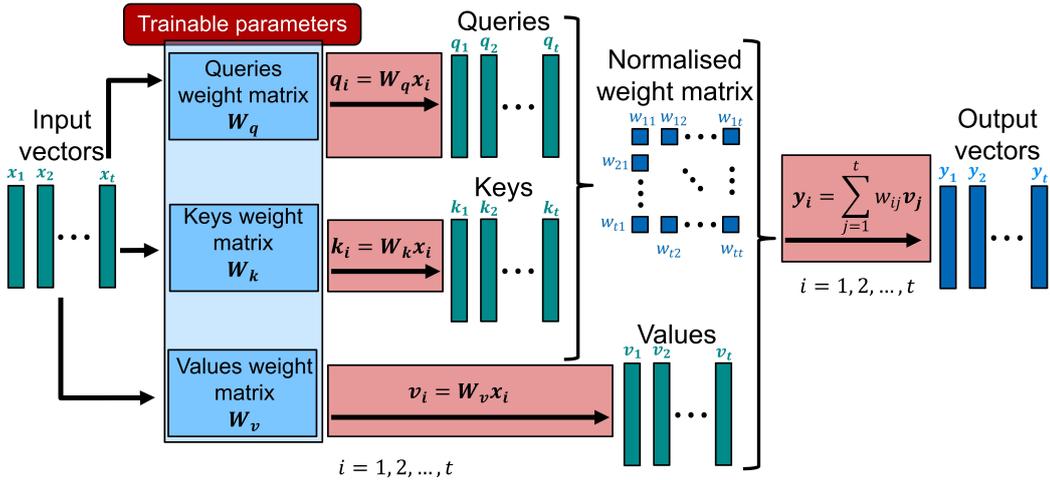


Fig. 6: Classic self-attention used in the classic Transformer (Vaswani et al., 2017). The input is t vectors of size K , while the generated output is a different set of t vectors of the same size. The input vectors together with the query, key, and value weight matrices are used to calculate queries, keys, and values; a weight matrix is produced by multiplying queries and keys, and it is subsequently scaled and normalised.

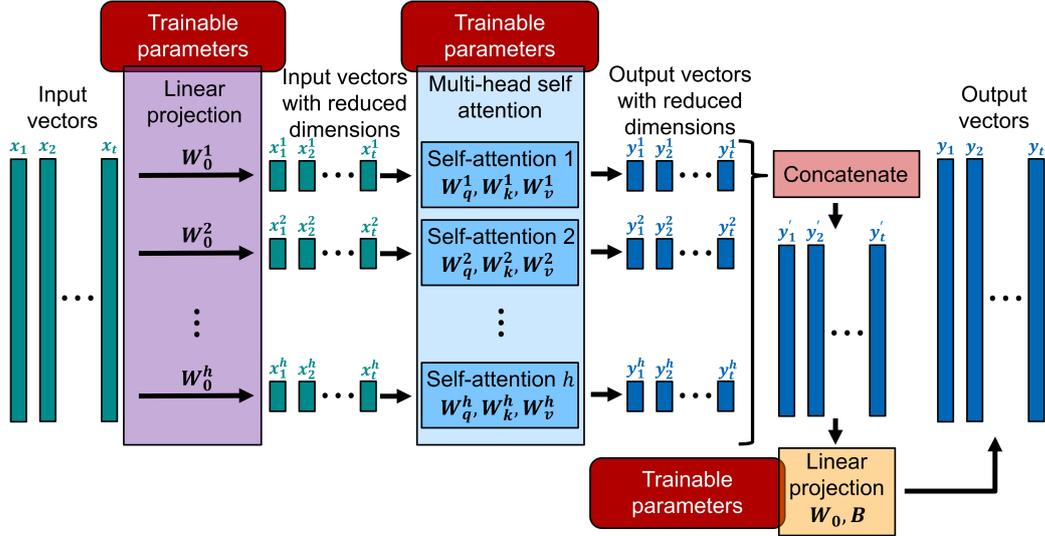


Fig. 7: Classic multi-head (self-)attention operation. Several self-attention operations are employed in parallel, which allows for a more efficient extraction of the various features in the given time series. The output vectors with reduced dimensions are concatenated and then the linear projection is applied to them to obtain the output vectors.

3 Methodology

3.1 Selected models

In this paper LSTM model and four Transformer-based models are applied to the 1D and 2D heat conduction problems. Table 1 provides a summary of these models; the classic Transformer will hereafter be referred to as the Transformer. The self-attention operation used in the Transformer (Vaswani et al., 2017) is described in Sub-section 2.4. Informer (Zhou et al., 2021), Autoformer (Wu et al., 2021), and FEDformer (Zhou et al., 2022) were developed with an aim of increasing the Transformer’s efficiency by reducing its complexity and adapting the architecture specifically to time-series processing (Table 1).

Table 1: The summary of the models considered in this paper.

Model type	LSTM	Transformer	Informer	Autoformer	FEDformer
Original purpose	Sequential data	Linguistic data	Temporal data	Temporal data	Temporal data
Self-attention type	N/A	Classic self-attention	Sparse self-attention	Auto-Correlation	Discrete Fourier Transform (DFT)

3.2 Model structure

Table 2 presents the model hyperparameters used for the Transformer-based models in this paper. Three layers, two encoder layers and one decoder layer, are used; three values of the input and output time-series length (sequence length hyperparameter) are considered: 25, 50, and 100. For greater clarity, the sequence length l will be referred to as the prediction window size, with the prediction window being defined as the time interval for which a prediction is made by the model. The hyperparameters No. 4-9 are assigned the values used in literature (Vaswani et al., 2017; Zhou et al., 2021; Wu et al., 2021; Zhou et al., 2022).

Table 3 presents the model hyperparameters used for the LSTM model in this paper. Only one layer is used; the model dimension is fixed at 512 to match the Transformer-based models. Three values of the prediction window size l are considered: 25, 50, and 100. Furthermore, the feed-forward layer is added after one LSTM layer in order to reshape an output and directly predict the temperature over multiple time steps in one inference step.

Table 2: Hyperparameters selected for all Transformer-based models.

No.	Hyperparameter	Value(s)
1.	Encoder layers	2
2.	Decoder layers	1
3.	Prediction window size (sequence length) l	25, 50, and 100
4.	Model dimension d_{model}	512
5.	Multi-head (self-)attention heads h	8
6.	Feed-forward network dimension	2048
7.	Dropout rate	0.05
8.	Activation function	GELU
9.	Attention factor	3

Table 3: Hyperparameters selected for LSTM models.

No.	Hyperparameter	Value(s)
1.	LSTM layers	1
2.	Prediction window size (sequence length) l	25, 50, and 100
3.	Model dimension d_{model}	512
4.	Dropout rate	0.05

3.3 Training

All Transformer-based models are trained using an Adam optimiser (Kingma and Ba, 2017). Additionally, the warm-up stage is applied to the learning rate as it was shown to improve the training process of the Transformer-based architectures (Xiong et al., 2020). The learning rate during the warm-up stage is given by:

$$lr(it) = \frac{it}{T_{warmup}} lr_{max} \quad \text{for } it \leq T_{warmup} \quad (8)$$

where lr is the learning rate, and it is an iteration number. While the learning rate after the warm-up stage is given by:

$$lr(it) = \frac{lr(it-1)\sqrt{T_{warmup}}}{\sqrt{it}} \quad \text{for } it > T_{warmup} \quad (9)$$

Transformer-based models are trained using lr_{max} and T_{warmup} values provided in Table 4 for 500 epochs with the batch size of 32; then the best option is selected for each model type listed in Table 1 using Normalised Root Mean Square Error (NRMSE) as a metric:

$$NRMSE = \frac{RMSE}{u_{true.mean}} \quad \text{and} \quad RMSE = \sqrt{\frac{\sum_{k=1}^{N_{total}} (u_{true.k} - u_{pred.k})^2}{N_{total}}} \quad (10)$$

where N_{total} is a total number of temperature data points, $u_{true.mean}$ is the true temperature mean, $u_{true.k}$ is the k^{th} true temperature value out of N_{total} temperature data points, while $u_{pred.k}$ is the k^{th} temperature value predicted by the ML model. Further details are provided in Section 4.

The LSTM model is trained using just an Adam optimiser (Kingma and Ba, 2017) for 500 epochs. Finally, NVIDIA A100 40GB GPU is used for training of all ML models considered in this paper. Figure 8 shows the convergence during the training, with the best options shown for Transformer, Informer, Autoformer, and FEDformer.

All models are initialised with a fixed random seed to ensure the repeatability of the results. This approach guarantees the consistency of the initial weights and any stochastic processes within

Table 4: Learning rates used to train the models for 500 epochs; the best option is selected for each model type listed in Table 1 using Normalised Root Mean Square Error (NRMSE) as a metric (Eq. 10).

Option No.	Constant lr or with warm-up	lr_{max}	T_{warmup}
1.	Constant $lr = 1e^{-4}$	N/A	N/A
2.	With warm-up	$1e^{-3}$	4000
3.	With warm-up	$1e^{-3}$	2000
4.	With warm-up	$1e^{-3}$	500
5.	With warm-up	$5e^{-4}$	4000
6.	With warm-up	$5e^{-4}$	2000
7.	With warm-up	$5e^{-4}$	500

the training algorithm across different runs. Thus, the potential variability in performance due to random initialization is mitigated allowing for a balanced comparison of the model architectures.

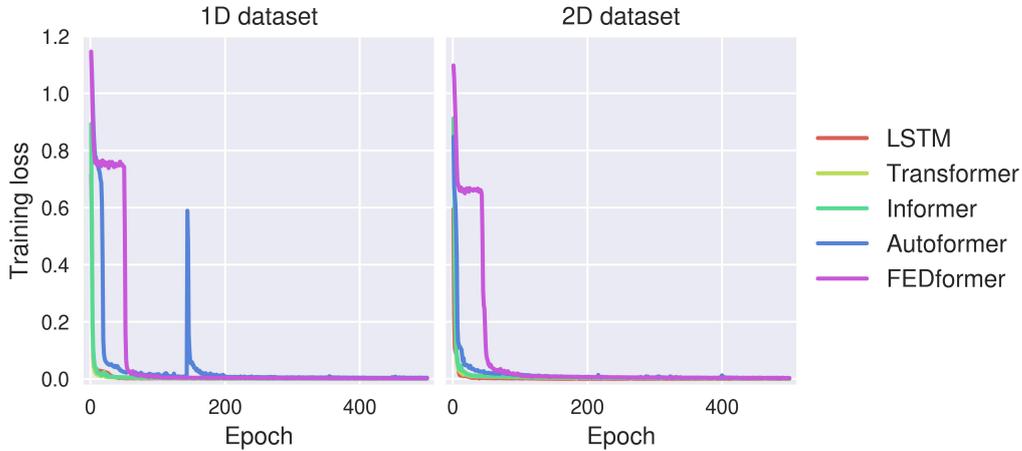


Fig. 8: Model convergence during the training process for 1D and 2D transient heat conduction cases discussed in Section 4.

4 Results and discussion

4.1 One-dimensional transient heat conduction

The linear 1D transient heat conduction equation is given by the following expression:

$$\frac{\partial u(x, t)}{\partial t} = \alpha \frac{\partial^2 u(x, t)}{\partial x^2} \quad (11)$$

where u is temperature, α is thermal diffusivity, t is time, and x is a space coordinate.

For this case α is set to be equal to $8.6e^{-4}m^2/s$, while the boundary and initial conditions are given by (Figure 9):

$$\begin{aligned} \text{Boundary Conditions: } & u(x = x_A, t) = 255.372K \quad \text{and} \quad \frac{\partial u(x = x_B, t)}{\partial x} = 0 \\ \text{Initial Conditions: } & u(x, t = 0) = 272.039K \end{aligned} \quad (12)$$

The Ground Truth (GT) solution is obtained using the Finite Difference Method (FDM) implemented in the PyPDE Python package (Zwicker, 2020). Figure 9 shows the grid used to generate the GT; Figure 10 shows the GT used for model testing. The simulation is run for 1000s and the temperature is recorded every second: the first 700s of the obtained data is used for ML training, the

next 100s for ML validation, and finally the last 200s is used for ML testing. The temperature values at the six input channels (Figure 9) are given to the ML model, whereas the temperature values at the 194 output channels are generated by the ML model (Figure 11). The locations of the six input channels are randomly selected.

For the training, validation, and testing the prediction window is moved by one time step forward, which is equal to one second in this case, to generate one input-output sample. For example, assuming that the prediction window size l is equal to 50, the first prediction window spans from 1s to 50s of 200s used for ML testing, the second one spans from 2s to 51s, the third from 3s to 52s etc. Consequently, for testing when $l = 50$ there are $200 - 50 + 1 = 151$ prediction windows for which the predicted temperature matrix $[U_{pred}]$ and true temperature matrix $[U_{true}]$ are generated. $[U_{pred}]$ and $[U_{true}]$ are matrices of dimension $((50 \cdot 151) \times N_{out})$, where N_{out} is the number of output channels, which is equal to 194 for this case. Therefore, for Eq. 10 N_{total} can be calculated as:

$$N_{total} = (50 \cdot 151) \cdot N_{out} = (50 \cdot 151) \cdot 194 \quad (13)$$

For other values of the prediction window size l and N_{out} the calculations are performed in the same manner.



Fig. 9: The 1D grid used to generate the Ground Truth (GT) for 1D transient heat conduction. The input and output channels of the mesh nodes are highlighted, with the six green crosses being the input channels and the 194 red dots being the output channels. The temperature values at the six input channels are given to the ML model, whereas the temperature values at the 194 output channels are generated by the ML model at every time step, which is equal to 1s in this case. The locations of the six input channels are randomly selected.

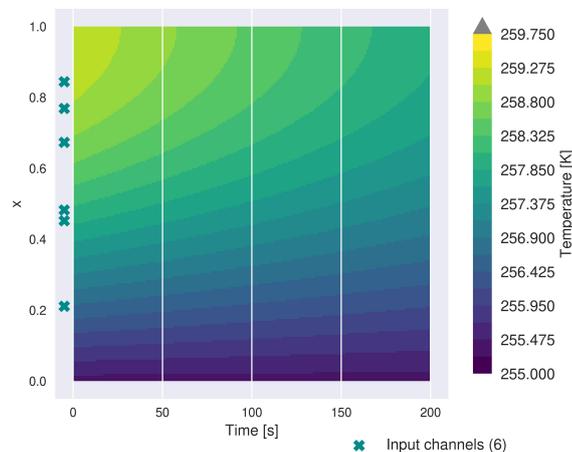


Fig. 10: The Ground Truth (GT) for 1D transient heat conduction problem generated using the Finite Difference Method (FDM) implemented in the PyPDE Python package (Zwicker, 2020). This region is used for the model testing. The temperature values at the six input channels (green crosses) are given to the ML model, whereas the temperature values at the 194 output channels (not shown on this figure, please refer to Figure 9) are generated by the ML model. The locations of the six input channels (green crosses) are randomly selected.

Table 5 provides testing NRMSEs calculated using Eq. 10 with Eq. 13, as well as the training times. In order to visualise the error distribution in space and time, four consecutive prediction windows (out of 151 prediction windows) are selected for the prediction window size $l = 50$; these are four prediction windows spanning from 1s to 50s of 200s used for ML testing, from 51s to 100s,

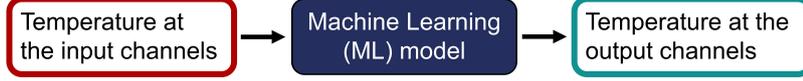


Fig. 11: Outline of the Machine Learning (ML) model used for the transient thermal field reconstruction. The temperature values at the six input channels are given to the ML model, whereas the temperature values at the 194 output channels are generated by the ML model. The locations of the six input channels are randomly selected.

from 101s to 150s, and finally from 151s to 200s. Figure 12 (top row) shows the prediction error distribution defined using Eq. 14, whereas the bottom row of this Figure shows the prediction errors averaged at each time step using Eq. 15.

$$PE_{i,n} = \frac{|u_{pred-i,n} - u_{true-i,n}|}{u_{true-i,n}} \quad (14)$$

$$PE_i = \frac{\sum_{n=1}^{N_{out}} PE_{i,n}}{N_{out}} \quad (15)$$

where $PE_{i,n}$ is a prediction error at node n at time step i , PE_i is a prediction error averaged at time step i ; $u_{pred-i,n}$ and $u_{true-i,n}$ are the predicted and true temperatures, respectively, at node n at time step i .

Table 5 shows that for all three values of l LSTM model achieves the lowest training times, whereas FEDformer model consistently has the highest training times. For $l = 25$ and $l = 50$ LSTM model attains the lowest NRMSE; and, indeed, this correlates with the prediction error distribution on Figure 12. For $l = 100$ Transformer displayed the lowest NRMSE; however, the NRMSE of LSTM is only 0.004% higher. Finally, Autoformer model obtained the NRMSE more than two times higher than LSTM model.

Table 5: Model testing errors and training times for the 1D heat conduction calculated using Eq. 10. The best results are highlighted in **green bold** and the worst results are highlighted with a **red underline**.

Model type	LSTM	Transformer	Informer	Autoformer	FEDformer
Prediction window size $l = 25$					
NRMSE (Eq. 10) [%]	0.188	0.189	0.192	<u>0.497</u>	0.362
Training time [min]	35.9	39.1	40.2	40.9	<u>46.4</u>
Prediction window size $l = 50$					
NRMSE (Eq. 10) [%]	0.189	0.190	0.190	<u>0.493</u>	0.365
Training time [min]	36.0	38.5	39.8	42.4	<u>52.3</u>
Prediction window size $l = 100$					
NRMSE (Eq. 10) [%]	0.196	0.192	0.195	<u>0.468</u>	0.371
Training time [min]	36.3	38.5	40.9	43.8	<u>65.2</u>

4.2 Two-dimensional transient heat conduction

The linear 2D transient heat conduction equation is given by the following expression:

$$\frac{\partial u(x, y, t)}{\partial t} = \alpha_x \frac{\partial^2 u(x, y, t)}{\partial x^2} + \alpha_y \frac{\partial^2 u(x, y, t)}{\partial y^2} \quad (16)$$

where x and y are the spatial coordinates, and α_x and α_y are the thermal diffusivities in x and y directions.

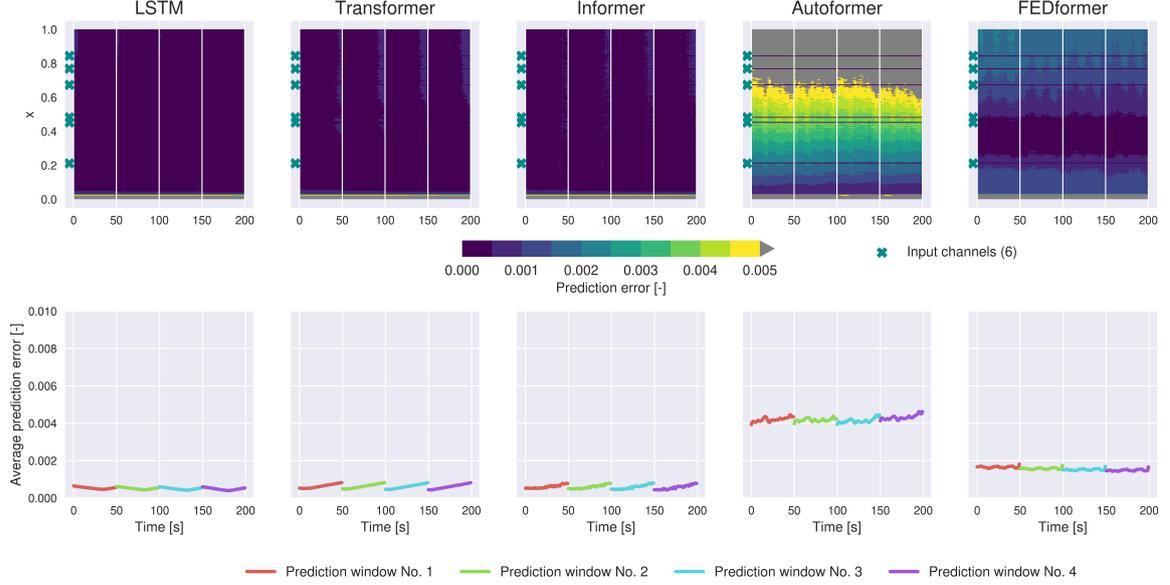


Fig. 12: Prediction error distribution calculated using Eq. 14 (top row) and prediction errors averaged at each time step calculated using Eq. 15 (bottom row) for the 1D heat conduction for five models with the prediction window size $l = 50$. For this error visualisation four prediction windows located consecutively to one another are selected. The six green crosses are the input channels.

For this case α_x is set to be equal to $13.9e^{-4}m^2/s$, while α_y is set to be equal to $3.3e^{-4}m^2/s$. The boundary and initial conditions are given by Eq. 17, and Figure 13 shows the location of the domain boundaries.

Boundary Conditions:

1. $u(x, y, t) = 255.372K$ for $x, y \in [AB] \cup [BC]$
2. $\nabla u(x, y, t) = 0$ for $x, y \in [CD] \cup [DA]$

$$\text{Initial Conditions: } u(x, y, t = 0) = 272.039K \quad (17)$$

The Ground Truth (GT) solution is obtained using the Finite Element Method (FEM) implemented in Code_Aster open-source software (Électricité de France (EDF), 1989–2023). Figure 13 shows the mesh used to generate the GT; Figure 14 shows the GT used for model testing. The simulation is run for 1000s and the temperature is recorded every second: the first 700s of the obtained data is used for ML training, the next 100s for ML validation, and finally the last 200s is used for ML testing. The temperature values at the twelve input channels are given to the ML model, whereas the temperature values at the 1142 output channels are generated by the ML model (Figures 11 and 13). The locations of the twelve input channels are randomly selected. Similar to 1D case (Sub-section 4.1), for the training, validation, and testing the prediction window is moved by one time step forward, which is one second in this case, to generate one input-output sample; the procedures for calculating the total number of prediction windows and NRMSEs are exactly the same as in the 1D case.

Table 6 provides testing NRMSEs calculated using Eq. 10 with Eq. 13, as well as the training times. The prediction error distributions for the prediction window size $l = 50$ are visualised in the similar way to 1D case (Sub-section 4.1). Figure 15 (top row) shows the time-averaged prediction error distribution defined using Eq. 18, whereas the bottom row of this figure shows the prediction errors averaged at each time step using Eq. 15 (please note that the different scales are used for these charts in comparison with Figure 12).

$$PE_n = \frac{\sum_{i=1}^{200} PE_{i,n}}{200} \quad (18)$$

where PE_n is a prediction error averaged at node n . Figure 16 shows the prediction error distribution variation with time; the four selected time instances correspond to the middle of each prediction window in Figure 15. This Figure highlights the areas where the prediction errors tend to increase with time.

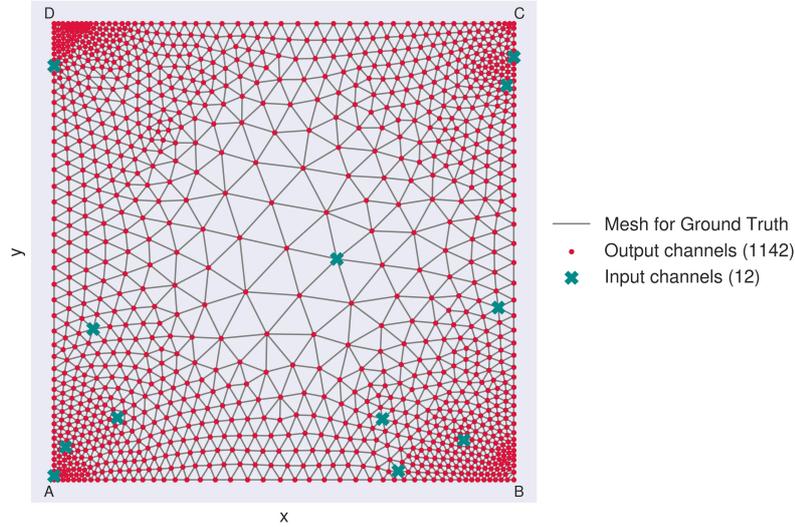


Fig. 13: The 2D mesh used to generate the Ground Truth (GT) for 2D transient heat conduction. The input and output channels of the mesh nodes are highlighted, with the twelve green crosses being the input channels and the 1142 red dots being the output channels. The temperature values at the twelve input channels are given to the ML model, whereas the temperature values at the 1142 output channels are generated by the ML model. The locations of the twelve input channels are randomly selected.

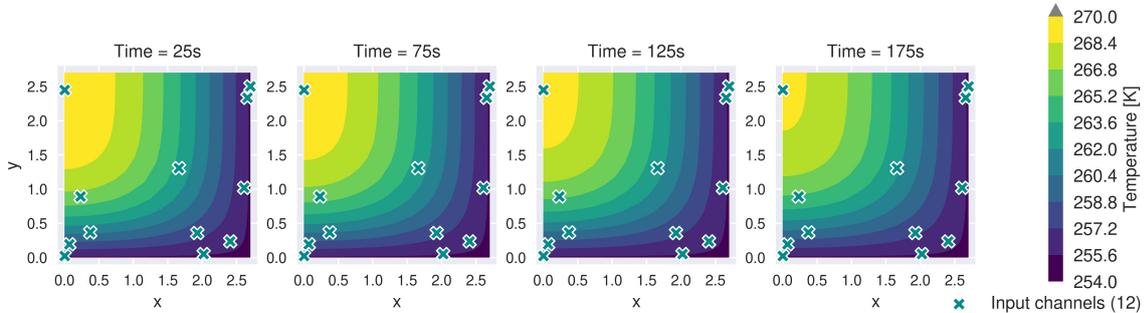


Fig. 14: The Ground Truth (GT) for 2D transient heat conduction problem generated using the Finite Element Method (FEM) implemented in Code_Aster open-source software (Électricité de France (EDF), 1989–2023). It is used for the model testing. The temperature values at the twelve input channels (green crosses) are given to the ML model, whereas the temperature values at the 1142 output channels (not shown on this figure, please refer to Figure 13) are generated by the ML model. The locations of the twelve input channels (green crosses) are randomly selected.

Table 6 shows that for all three values of l FEDformer displayed the lowest NRMSEs and the highest training times. However, the NRMSEs achieved by LSTM model is only 0.1-0.2% higher than FEDformer, whilst LSTM model’s training is 20-41% faster than FEDformer’s. Moreover, one-layer LSTM architecture is significantly less complex than FEDformer structure, meaning that it is easier to troubleshoot it.

Surprisingly, the overall prediction error distribution patterns display similar features for all models (Figures 15 and 16). For the Transformer-based models this can be potentially explained by

the fact that all these models are structured around self-attention operation of some description. However, the fact that the simple LSTM model, the structure of which is unrelated to the self-attention operation, demonstrates almost identical error patterns is rather counterintuitive and may potentially challenge the confidence in Transformer-based models.

Table 6: Model testing errors and training times for the 2D heat conduction calculated using Eq. 10. The best results are highlighted in **green bold** and the worst results are highlighted with a **red underline**.

Model type	LSTM	Transformer	Informer	Autoformer	FEDformer
Prediction window size $l = 25$					
NRMSE (Eq. 10) [%]	2.197	2.208	2.218	<u>2.543</u>	2.015
Training time [min]	37.9	38.3	40.9	42.7	<u>47.5</u>
Prediction window size $l = 50$					
NRMSE (Eq. 10) [%]	2.170	2.191	2.193	<u>2.415</u>	2.015
Training time [min]	39.7	39.6	41.5	44.2	<u>67.0</u>
Prediction window size $l = 100$					
NRMSE (Eq. 10) [%]	2.127	2.155	2.157	<u>2.430</u>	1.979
Training time [min]	42.0	42.0	43.6	48.1	<u>66.9</u>

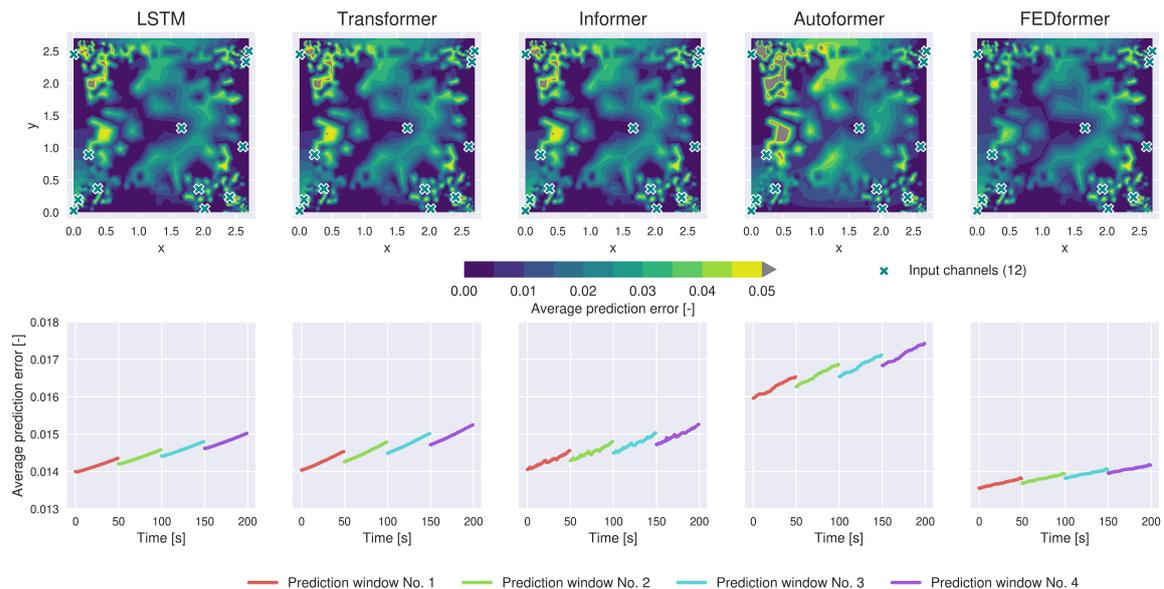


Fig. 15: Time-averaged prediction error distribution calculated using Eq 18 (top row) and prediction errors averaged at each time step calculated using Eq. 15 (bottom row) for the 2D heat conduction for five models with the prediction window size $l = 50$. For this error visualisation four prediction windows located consecutively to one another are selected. The twelve green crosses are the input channels.

5 Conclusions

In conclusion, the popular Transformer-based ML models are compared with the simple one-layer LSTM models for transient thermal field reconstruction problems. Generally, there are three main reasons why it might be advantageous to use the Transformer-based models over RRNs, such as LSTM, for sequence processing (Vaswani et al., 2017):

1. Self-attention operation is more effective at capturing the long-range dependencies in a sequence, thus the Transformer-based models should be more accurate.
2. They are more parallelisable as there are no recurrent connections; consequently, the training time should be shorter.
3. Due to the attention maps (Appendix A), they are more interpretable, and thus suffer less from the “black box” syndrome (Rudin, 2019), which ML models tend to be afflicted with; this should facilitate better trust in the predictions.

However, the present work demonstrates that none of the aforementioned reasons hold true for the considered cases of thermal problems concerned with solution reconstruction by providing the following arguments:

1. Transformer-based models exhibited lower or comparable prediction accuracy relative to the simple LSTM model.
2. The training times of the Transformer-based models are higher or comparable to the simple LSTM model.
3. The attention maps (Appendix A) may provide interpretability improvement for the language-related tasks, such as NLP (Vaswani et al., 2017) and computer vision (Kolesnikov et al., 2021), since the significance of the attention weights in attention maps can be easily intuitively deduced from the relationship between words in a sentence and from the parts of an image the attention operation “pays attention” to, respectively. However, this is not so easily done when dealing with just temporal data, as the relationship between values at different time steps cannot be interpreted in such an intuitive fashion, especially for the problems intrinsically based on a set of differential equations. Moreover, the LSTM layer weights can be visualised in the similar manner (Appendix B) and have the same rather low level of interpretability for the problems considered in this paper. Consequently, the interpretability advantage is negated.

Overall, the results of this study suggest that there is no significant benefit to using complex Transformer-based ML models over the conventional simpler ML models, such as classic LSTM network, for solving transient thermal field reconstruction problems.

Regarding the application of these models to the practical problems, such as the one outlined in Sub-section 2.1, the appropriate number of the reliable forward simulations would need to be available in order to generate the training data. Or, alternatively, the abundance of experimental data would need to be collected, which is rarely possible. Additionally, the computational effort of the training process and the memory usage increase with the length of the input time sequence (Table 7); consequently, for more complex problems significantly more time should be allocated for the training and VRAM (GPU memory) more carefully managed with batching.

Table 7: Computational complexity and memory usage of the models considered in this paper.

Model type	LSTM	Transformer	Informer	Autoformer	FEDformer
Computation complexity (training)	$O(L)$	$O(L^2)$	$O(L \log L)$	$O(L \log L)$	$O(L)$
Memory usage (training)	$O(L)$	$O(L^2)$	$O(L \log L)$	$O(L \log L)$	$O(L)$

* L represents the length of the input time sequence.

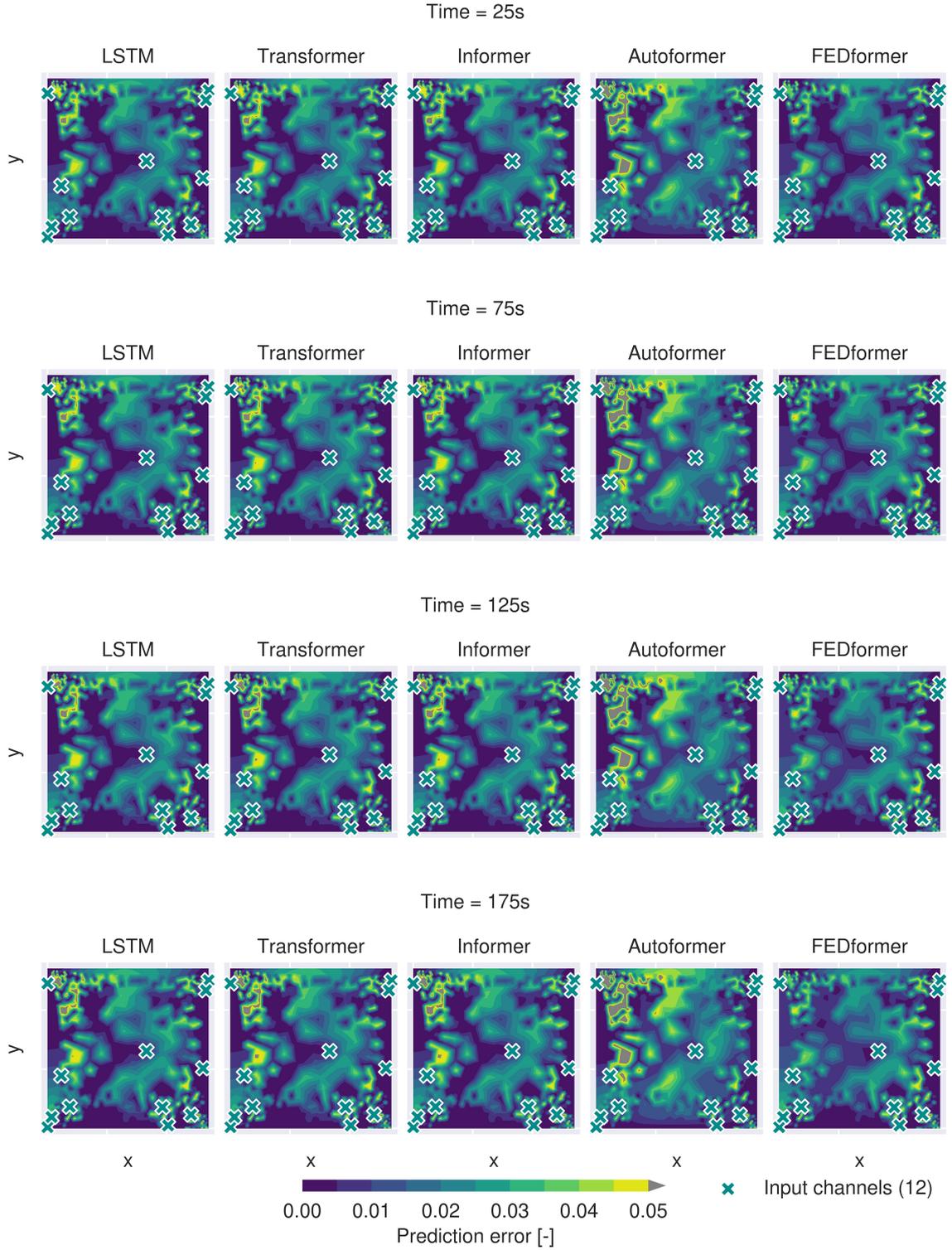


Fig. 16: Prediction error distribution dependence on time for the 2D heat conduction for five models with the prediction window size $l = 50$. The four selected time instances correspond to the middle of each prediction window in Figure 15. The prediction errors are calculated using Eq. 14 for fixed time instances i .

6 Acknowledgements

This work was part-funded by the United Kingdom Atomic Energy Authority (UKAEA) and the Engineering and Physical Sciences Research Council (EPSRC) under Grant Agreement Numbers

EP/W006839/1, EP/T517987/1 and EP/R012091/1. We acknowledge the support of Supercomputing Wales and AccelerateAI projects, which is part-funded by the European Regional Development Fund (ERDF) via the Welsh Government for giving us access to NVIDIA A100 40GB GPUs for batch training.

References

- Alché, F. and de La Fortelle, A., An LSTM network for highway trajectory prediction, *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 353–359, 2017.
- Arridge, S., Maass, P., Öktem, O., and Schönlieb, C.-B., Solving inverse problems using data-driven models, *Acta Numerica*, vol. 28, pp. 1–174, 2019.
- Bangian-Tabrizi, A. and Jaluria, Y., An optimization strategy for the inverse solution of a convection heat transfer problem, *International Journal of Heat and Mass Transfer*, vol. 124, p. 1147 – 1155, 2018.
- Bengio, Y., Simard, P., and Frasconi, P., Learning long-term dependencies with gradient descent is difficult, *IEEE Transactions on Neural Networks*, vol. 5, no. 2, pp. 157–166, 1994.
- Bloem, P., Transformers from scratch, 2019, URL <https://peterbloem.nl/blog/transformers>, accessed on: March 29, 2023.
- Électricité de France (EDF), Code_Aster: Open Source Finite Element Solver, Analysis of Structures and Thermomechanics for Studies and Research, Open source on <https://code-aster.org/>, 1989–2023.
- Gers, F. A., Schmidhuber, J., and Cummins, F., Learning to Forget: Continual Prediction with LSTM, *Neural Computation*, vol. 12, no. 10, pp. 2451–2471, 2000.
- Goodfellow, I., Bengio, Y., and Courville, A., *Deep Learning*, The MIT Press, 2016.
- Hancock, D., Homfray, D. A., Porton, M., Todd, I., Wynne, B. P., Bamber, R., Flinders, K., Jepson, P., and Lewtas, H., Testing Advanced Divertor Concepts for Fusion Power Plants Using a Small High Heat Flux Facility, *Journal of Testing and Evaluation (ASTM International)*, 2018.
- Hochreiter, S. and Schmidhuber, J., Long Short-term Memory, *Neural computation*, vol. 9, pp. 1735–80, 1997.
- Hsu, W.-N., Zhang, Y., Lee, A., and Glass, J. R., Exploiting Depth and Highway Connections in Convolutional Recurrent Deep Neural Networks for Speech Recognition, in *Interspeech*, 2016.
- Ivanovic, B. and Pavone, M., The Trajectron: Probabilistic Multi-Agent Trajectory Modeling With Dynamic Spatiotemporal Graphs, in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- Kingma, D. P. and Ba, J., Adam: A Method for Stochastic Optimization, 2017.
- Kolesnikov, A., Dosovitskiy, A., Weissenborn, D., Heigold, G., Uszkoreit, J., Beyer, L., Minderer, M., Dehghani, M., Houlsby, N., Gelly, S., Unterthiner, T., and Zhai, X., An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, 2021.
- Lim, B. and Zohren, S., Time-series forecasting with deep learning: a survey, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 379, no. 2194, p. 20200209, 2021.
- Qin, C., Zhang, A., Zhang, Z., Chen, J., Yasunaga, M., and Yang, D., Is ChatGPT a General-Purpose Natural Language Processing Task Solver?, 2023.

- Ren, L., Dong, J., Wang, X., Meng, Z., Zhao, L., and Deen, M. J., A Data-Driven Auto-CNN-LSTM Prediction Model for Lithium-Ion Battery Remaining Useful Life, *IEEE Transactions on Industrial Informatics*, vol. 17, no. 5, pp. 3478–3487, 2021.
- Rudin, C., Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead, *Nature Machine Intelligence*, vol. 1, pp. 206–215, 2019.
- Tamaddon-Jahromi, H. R., Chakshu, N. K., Sazonov, I., Evans, L. M., Thomas, H., and Nithiarasu, P., Data-driven inverse modelling through neural network (deep learning) and computational heat transfer, *Computer Methods in Applied Mechanics and Engineering*, vol. 369, p. 113217, 2020.
- Tarantola, A., *Inverse Problem Theory and Methods for Model Parameter Estimation*, Society for Industrial and Applied Mathematics, 2004.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I., Attention is All you Need, in *Advances in Neural Information Processing Systems*, Curran Associates, Inc., 2017, vol. 30.
- Wen, Q., Zhou, T., Zhang, C., Chen, W., Ma, Z., Yan, J., and Sun, L., Transformers in Time Series: A Survey, 2023.
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Le Scao, T., Gugger, S., Drame, M., Lhoest, Q., and Rush, A., Transformers: State-of-the-Art Natural Language Processing, in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Online: Association for Computational Linguistics, 2020, pp. 38–45.
- Wu, B., Wang, L., and Zeng, Y.-R., Interpretable wind speed prediction with multivariate time series and temporal fusion transformers, *Energy*, vol. 252, p. 123990, 2022.
- Wu, H., Xu, J., Wang, J., and Long, M., Autoformer: Decomposition Transformers with Auto-Correlation for Long-Term Series Forecasting, in *Neural Information Processing Systems*, 2021.
- Xiong, R., Yang, Y., He, D., Zheng, K., Zheng, S., Xing, C., Zhang, H., Lan, Y., Wang, L., and Liu, T., On Layer Normalization in the Transformer Architecture, in H. D. III and A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, PMLR, 2020, vol. 119 of *Proceedings of Machine Learning Research*, pp. 10524–10533.
- Yu, Y., Si, X., Hu, C., and Zhang, J., A Review of Recurrent Neural Networks: LSTM Cells and Network Architectures, *Neural Computation*, vol. 31, no. 7, pp. 1235–1270, 2019.
- Zhang, Y., Wang, S., Ji, G. et al., A comprehensive survey on particle swarm optimization algorithm and its applications, *Mathematical problems in engineering*, vol. 2015, 2015.
- Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., and Zhang, W., Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting, *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 12, pp. 11106–11115, 2021.
- Zhou, T., Ma, Z., Wen, Q., Wang, X., Sun, L., and Jin, R., FEDformer: Frequency Enhanced Decomposed Transformer for Long-term Series Forecasting, in K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, PMLR, 2022, vol. 162 of *Proceedings of Machine Learning Research*, pp. 27268–27286.
- Zwicker, D., py-pde: A Python package for solving partial differential equations, *Journal of Open Source Software*, vol. 5, p. 2158, 2020.

Appendix A Attention maps for Transformer-based models

In general, self-attention operation or a variation thereof allows the model to attend to different parts of the input vector sequence in order to generate an output vector sequence, meaning that it concentrates more on the selected input vectors, which it deems to be more relevant, while generating a certain output vector (Vaswani et al., 2017). As it is mentioned in Sub-section 2.4, the queries and keys produced using query and key weight matrices are combined to produce a weight matrix, which can be called an attention matrix; each element of this attention matrix, which can be referred to as an attention score, represents the contribution it makes towards a certain output vector.

In order to illustrate this, the NLP example shown in Figure A1 can be considered (Bloem, 2019). The key representing the qualities the book contains and the query, which represent the reader preferences, are matched using a dot product to obtain the attention score. This score demonstrates how well the book matches the reader’s preferences. Generally, the attention score indicates the degree of relevance between the key and the query; thus, it shows to what degree a certain output vector out of the output vector sequence is influenced by a certain input vector out of the input vector sequence. In the case of the transient thermal field reconstruction problems considered in this paper, each input vector i contains the information provided by the input channels at i^{th} time step, while each output vector i contains the information provided by the output channels at i^{th} time step (Fig. 9 and 13). For the classic Transformer, the attention score indicating the degree to which the output vector i is influenced by the input vector j can be calculated using a dot product between the query i and the key j (Eq. A1).

$$Attention_score_{(i,j)} = \{\mathbf{q}_i\}^T \cdot \{\mathbf{k}_j\} \quad i, j = 1, 2, \dots, t \quad (\text{A1})$$

where $\{\mathbf{q}_i\}$ is the query vector i corresponding to the output vector i , $\{\mathbf{k}_j\}$ is the key vector j corresponding to the input vector j , and t is a sequence length (Table 2).

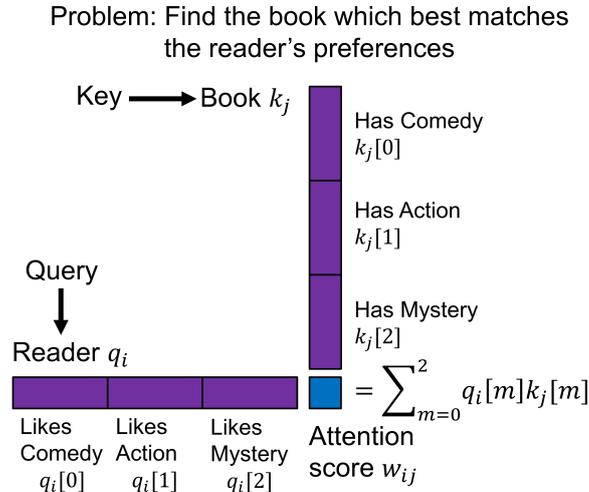


Fig. A1: The book example illustrating the self-attention operation. The key representing the qualities the book contains and the query, which represent the reader preferences, are matched using a dot product to obtain the attention score. This score demonstrates how well the book matches the reader’s preferences. Generally, the attention score indicates the degree of relevance between the key and the query; thus, it shows to what degree a certain output vector out of the output vector sequence is influenced by a certain input vector out of the input vector sequence. In the case of the transient thermal field reconstruction problems considered in this paper, each input vector i contains the information provided by the input channels at i^{th} time step, while each output vector i contains the information provided by the output channels at i^{th} time step (Figures 9 and 13).

The attention matrix containing the attention scores can be visualised as a map by projecting it on an image where each cell (i, j) corresponds to the attention score computed for query vector i and key vector j . Figure A2 shows an example of an attention score (i, j) located on the attention map image.

Transformer-based models considered in this paper employ a multi-head attention operation (Figure 7) where several self-attention operations are performed in parallel (Sub-section 2.4). Consequently, each multi-head attention operation generates a number of attention maps equal to the number of self-attention heads the multi-head attention operation is comprised of. Eight attention heads are used in this paper (Table 2); therefore, eight attention maps can be generated per one multi-head attention operation. It is important to note that these attention scores should not be used to compare the different heads between each other, they should only be used to compare the attention scores from one attention map between each other. Furthermore, the attention matrix size for Transformer (Vaswani et al., 2017) and Informer (Zhou et al., 2021) is prediction window size l by prediction window size l (Table 2); however, the attention matrix size for Autoformer (Wu et al., 2021) and FEDformer (Zhou et al., 2022) is sequence length by d_{model}/h (Table 2), this is because the keys are projected in the sequence length direction prior to the calculation of attention scores.

Figures A3 and A4 show four out of eight attention maps for multi-head attention in the first layer of each model for $l = 50$. The attention scores are scaled between 0 and 1 with the scaler fitting done separately for each attention map.

Some distinctive characteristics can be observed in these attention maps. Transformer and Informer attention maps are diverse, which means that different attention heads tend to attend to different parts of the input vector sequence. This is beneficial for the model’s performance because it signifies that the attention heads are extracting different features from the input data. Transformer’s attention maps are smoother than Informer’s ones, which can be attributed to the fact that Informer utilises a sparse version of Transformer’s self-attention operation. On the other hand, in spite of some differences Autoformer’s attention maps look vastly similar, which is undesirable since this means that the attention heads are attempting to extract the same features from the input data. The attention map similarity might be the reason why Autoformer’s performance is worse than the performance displayed by other models.

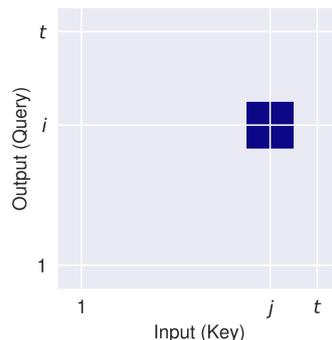


Fig. A2: Attention score example. The attention matrix containing the attention scores can be visualised as a map by projecting it on an image where each cell (i, j) corresponds to the attention score computed for query vector i and key vector j . This Figure shows an example of an attention score (i, j) located on the attention map image.

Overall, the ability to produce attention maps allows the Transformer-based model to be more interpretable than other ML models, as they show how the model focuses on the various sections of the input vector sequence. By visualising the the attention scores in this manner, it is possible to see which parts of the input sequence the model “pays attention” to at each time step.

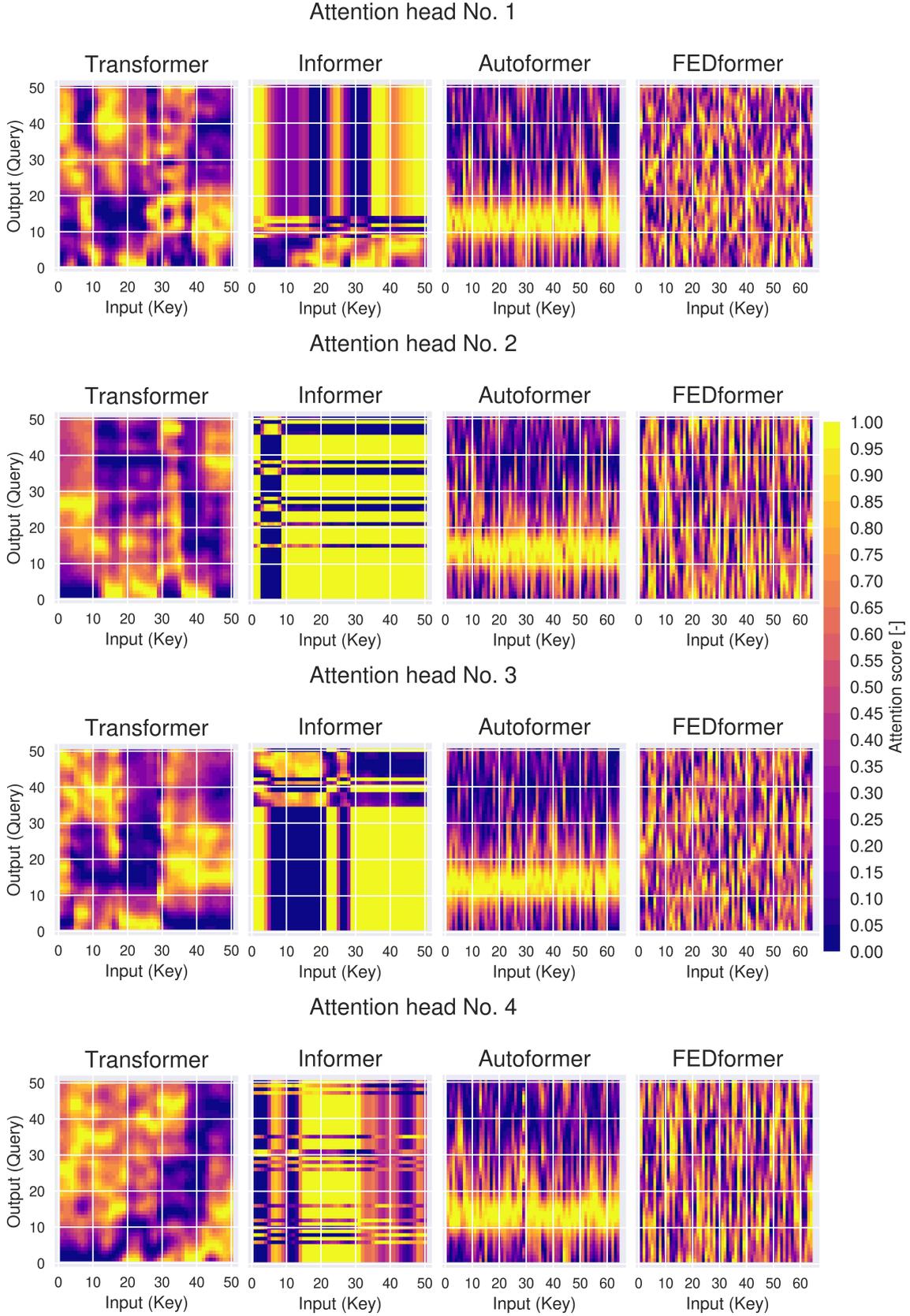


Fig. A3: Attention maps for Encoder layer No. 1 for 1D heat conduction problem for four models considered in this paper. The prediction window size l is equal to 50. This Figure show four out of eight attention maps for multi-head attention in the first layer of each model. The attention scores are scaled between 0 and 1 with the scaler fitting done separately for each attention map.

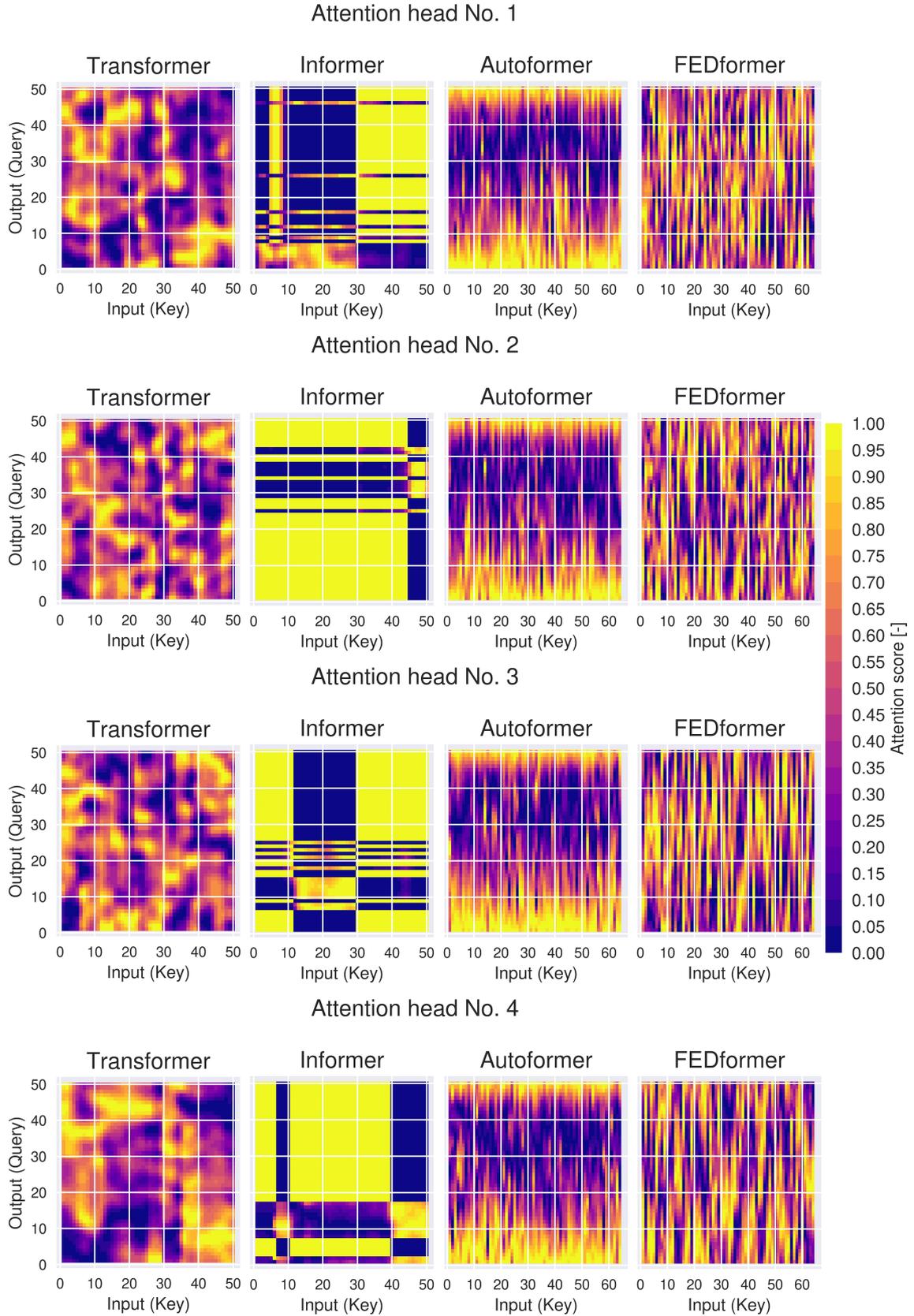


Fig. A4: Attention maps for Encoder layer No. 1 for 2D heat conduction problem for four models considered in this paper. The prediction window size l is equal to 50. This Figure shows four out of eight attention maps for multi-head attention in the first layer of each model. The attention scores are scaled between 0 and 1 with the scaler fitting done separately for each attention map.

Appendix B Weight visualisations for LSTM models

The purpose of this appendix is to illustrate that the weight visualisations similar to the attention maps can be easily constructed for the one-layer LSTM model. These LSTM weight visualisations have the same low level of intuitive interpretability as for the Transformer-based models specifically for the transient thermal field reconstruction problems considered in this paper. Figure B5 shows the values of h_t for each time step l (Eq. 3) for LSTM models with the prediction window size equal to 50; these values are scaled between 0 and 1 with the scaler fitting done separately for each visualisation. For consistency these values are referred to as attention scores in Figure B5.

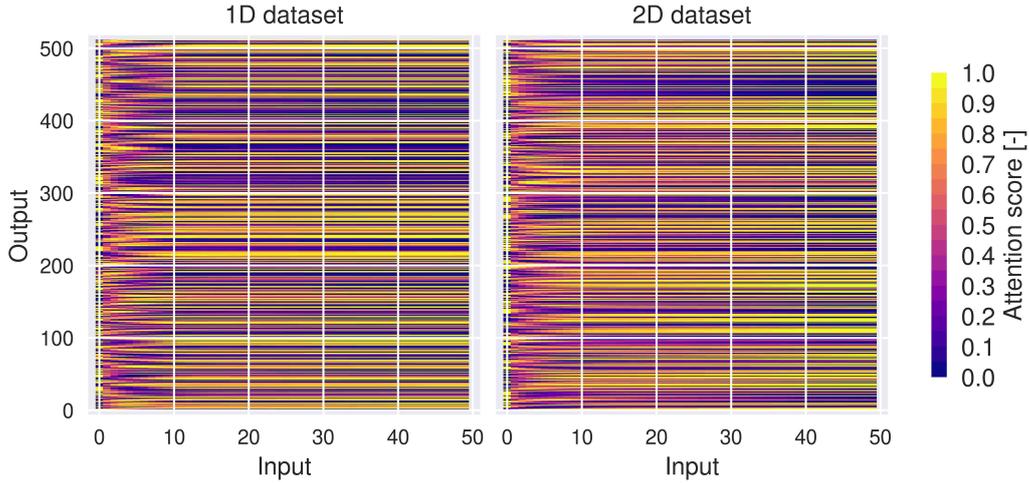


Fig. B5: Weight visualisations for 1D and 2D heat conduction problems for LSTM models with the prediction window size equal to 50. This Figure shows the values of h_t for each time step l (Eq. 3). The attention scores are scaled between 0 and 1 with the scaler fitting done separately for each weight visualisation.