Simultaneous Multistep Transformer Architecture for Model Predictive Control

Junho Park^a, Mohammad Reza Babaei^a, Samuel Arce Munoz^a, Ashwin N. Venkat^b, John D. Hedengren^a

^aDepartment of Chemical Engineering, Brigham Young University, Provo, Utah, USA ^bSeeq Corporation, Seattle, Washington, USA

Abstract

Transformer neural networks have revolutionized natural language processing by effectively addressing the vanishing gradient problem. This study focuses on applying Transformer models to time-series forecasting and customizing them for a simultaneous multistep-ahead prediction model in surrogate model predictive control (MPC). The proposed method showcases improved control performance and computational efficiency compared to LSTM-based MPC and one-step-ahead prediction models using both LSTM and Transformer networks. The study introduces three key contributions: (1) a new MPC system based on a Transformer time-series architecture, (2) a training method enabling multistep-ahead prediction for time-series machine learning models, and (3) validation of the enhanced time performance of multistep-ahead Transformer MPC compared to one-stepahead LSTM networks. Case studies demonstrate a significant fifteen-fold improvement in computational speed compared to one-step-ahead LSTM, although this improvement may vary depending on MPC factors like the lookback window and prediction horizon.

Keywords: Transformer neural network architecture, Long short-term memory, model predictive control

1 1. Introduction

Machine learning technologies have gained attention in many industries and have grown in popularity, especially in areas requiring classification, such as pattern recognition for image and voice. Many researchers in process system engineering have recently renewed interest in machine learning technologies, such as

Preprint submitted to Computers & Chemical Engineering

August 15, 2023

deep learning. An overview of current trends and opportunities for machine learn-6 ing technologies in process system engineering has been discussed in [1, 2, 3, 4]. 7 Areas using machine learning technologies in process system engineering are con-8 trol and optimization. Model predictive control (MPC) and real-time optimization 9 (RTO) are two advanced automation systems that require process models. Many 10 MPC application research studies have been conducted using artificial neural net-11 work (ANN) models: internal combustion engine [5], distillation column and con-12 tinuous stirred tank reactor (CSTR) example [6, 7], Heating ventilation and air 13 conditioning systems (HVAC) [8, 9], LSTM (Long short term memory network) 14 MPC for continuous pharmaceutical reactors [10], and phthalic anhydride syn-15 thesis in a fixed-bed catalytic reactor [11]. A Hybrid model with a Hammerstein 16 model structure, which consists of a static neural net model block and a linear 17 dynamics model block in series, has been investigated in [12]. A gated recurrent 18 unit (GRU)-based encoder-decoder architecture integrated with a physics-based 19 model has been tested for thermal power plants [13]. An approximated MPC con-20 cept has also been investigated by learning control policy from the closed-loop 21 performance data [14]. Reinforcement learning methods for RTO application have 22 been investigated in [15, 16, 17]. 23

Among the different types of ANNs, RNNs are a natural choice for model-24 ing time-series data because the structure of RNN states has explicit temporal 25 dependence. However, the sequential processing nature of the RNNs makes the 26 model prediction computationally expensive by recursively updating the hidden 27 state. A novel time-series ANN architecture, Transformer neural network, has 28 been introduced to solve the issue by parallelizing the model prediction steps, 20 making it possible to produce multistep predictions simultaneously [18]. State-30 of-the-art ANN architectures, especially attention mechanisms, improve natural 31 language processing (NLP). Although dynamic process models can be built with 32 time-series data like the NLP problems, the actual issues are vastly different in 33 model usage for prediction in model-based control and optimization applications. 34 Previous research has recognized the Transformer network for its remark-35 able computational speed. A distinguishing feature of this network is its atten-36 tion mechanism, which adeptly handles irregular temporal dependencies. This 37 capability is particularly notable in the context of natural language data. This 38 study further explores the potential for the application of the Transformer net-39 work within process control fields by integrating it into the MPC framework as a 40 predictive model. In addition to creating a Transformer-based MPC system, the 41 research introduces a multistep-ahead prediction structure to maximize the bene-42 fits derived from the Transformer advancements, particularly within the context of 43

MPC applications. Evidence from three case studies indicates that the multistep
 ahead Transformer MPC offers superior performance in terms of computational
 efficiency and setpoint tracking compared to the one-step ahead LSTM MPC.

In an effort to contextualize these findings, the study contrasts the proposed 47 multistep ahead Transformer MPC with the one-step ahead LSTM, which main-48 tains the traditional data-driven time-series format within the surrogate MPC frame-49 work. Although all the case studies do not exhaustively examine two alternative 50 model structures — multistep ahead LSTM and one-step ahead Transformer — 51 the first case study (5.1) provides an indirect comparison of these models to the 52 multistep ahead Transformer MPC in terms of prediction accuracy and computa-53 tional efficiency. Further, sequential neural network models such as LSTM are 54 recognized for a relative lack of prediction precision, particularly when dealing 55 with irregular temporal dependencies prevalent in the multistep-ahead data struc-56 ture as shown in Figure 7. 57

While this paper does not delve into a comprehensive investigation of all model structures, it lays the groundwork for future research. Further studies could provide valuable insights by examining the performance of various models in relation to complexity and ability to handle data irregularities. This could pave the way for a more understanding of model behaviors in different scenarios.

The objective of this work is to develop a simultaneous multistep-ahead pre-63 diction method using a Transformer architecture for MPC applications. Figure 64 1 gives an overview of two approaches for using Transformer models for model 65 based-control. The Transformer model can emulate the process model (Figure 1a) 66 or the entire model predictive control (Figure 1b). For the scope of this paper, 67 the exploration primarily focuses on the former, investigating the surrogate MPC 68 approach more intensively. A brief review of Transformer Network Architecture 69 is discussed with details of the adaptations for time-series data. Source code is 70 available from https://github.com/BYU-PRISM/Transformer_MPC 71

72 2. Transformer Network Architecture

The original concept of Transformer architecture in "Attention is all you need" [18] is dedicated to machine translation tasks (Figure 2). Thus, each element of the original design needs to be customized for the time-series data processing. This section discusses each part of the original Transformer design and the modifications for time-series data, especially for the simultaneous multistep-ahead prediction needed for MPC.



(b) Emulation control structure

Figure 1: Two scenarios for Transformers in model based control



Figure 2: Transformer structure

79 2.1. Encoder-Decoder

The Encoder-Decoder concept is initially used in the Sequence to Sequence (Seq2seq) architecture that places two separate RNNs in a row. The Encoder RNN processes the input sequence to produce a context vector. This vector, which encapsulates the information from the input sequence, is then passed to the Decoder RNN to predict the output sequence. The Transformer architecture retains the encoder-decoder framework while using new concepts for specific components.

86 2.2. Positional Encoding

In the RNN models, the hidden state of the previous timesteps returns to the 87 network as an additional input along with the actual input values for the next time 88 step. This recursive data processing is not computationally efficient, but it has 89 an advantage in that it can inherently recognize the order of the input sequence. 90 However, as the Transformer model takes the whole input sequence simultane-91 ously, the inherent temporal dependency in the RNN is no longer available in the 92 Transformer. A positional embedding layer explicitly labels the sequence position 93 to have the Transformer network recognize the sequence order. Trigonometric 94 functions [18] is a method among many to embed the positional information into 95

the input sequence. This method is particularly effective for sequences that have an irregular length and interval by assigning the unique index in the form of a real number vector. Thus, this method is commonly used for NLP tasks that vary the number of words in every sentence, resulting in different input vector lengths and time intervals. On the other hand, the time-series data usually has a fixed size and interval or can easily be converted to this form by signal processing.

102 2.3. Attention Mechanism

The attention mechanism is a crucial part of the Transformer model. It has 103 been proposed to improve the long-term memory capability in the Seq2seq model 104 [19]. The encoded input sequence in the Seq2seq model is saved in a fixed-length 105 vector called context vector and passed to the decoder RNN [20]. Because the 106 context vector in the Seq2seq model is just the last hidden state of the encoder 107 RNN, it cannot accommodate the most critical correlation at earlier time steps 108 caused by vanishing gradients. The memory issue of the RNN can be significant 109 when the time series is long. The attention mechanism improves the memory 110 of the context vector by introducing a probability distribution between a hidden 111 state of the decoder RNN at a specific time step and the hidden states of every 112 time step in the encoder RNN. This probability distribution helps the model focus 113 more on a certain time step in the input sequence regardless of the order. It is 114 converted further into a vector with the same size as the hidden state, called the 115 attention value. This attention value acts as a context vector of the Seq2seq model, 116 holding more useful long-term correlation information between the encoder and 117 decoder. The Transformer architecture extensively employs the attention mecha-118 nism, which removes the inherent sequential processing in the RNN models. A 119 conditional probability of a sequence data for the multistep-ahead prediction is 120 shown in Equation 1. 121

$$\mathbf{p}(\mathbf{y}_{k+1:k+P} | \mathbf{y}_{k-w:k}, \mathbf{u}_{k-w:k+P}) = \prod_{t=k+1}^{k+P} \mathbf{p}(\mathbf{y}_t | \mathbf{y}_{t-w-1:t-1}, \mathbf{u}_{t-w-1:t-1}),$$
(1)

where, $\mathbf{p}(A \mid B)$ represents the probability of *A* under the given condition *B*. **y** and u denote the system output and system input, respectively. The past system output y is used as one of the NN model inputs, called auto-regressive input, while the system input **u** is called exogenous input. The future values (t > k) of **u** are also included in the exogenous input as known values for the NN model training, while the **y** only include the past values (t < k) as shown in LHS of Equation 1. *k*, *w*,

and P denote the current time step, look-back window size, and prediction horizon 128 size, respectively. Using the RNN models to achieve the multistep prediction, it 129 is inevitable to build up the hidden state recursively from time step k - w to k and 130 repeat the same recursive computations for every prediction value in the predic-131 tion horizon P. However, the attention mechanism simultaneously computes the 132 probability distribution, also called 'Attention distribution' between the input and 133 output sequence, altering the recurrence in RNN with the matrix multiplications 134 and *softmax* activation function. 135

Attention
$$(Q, K, V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V,$$
 (2)

where, Q, K, and V represent the weighted sequence data matrices, Query, Key, 136 and Value. The same input sequence data, X, is assigned for all Q, K, and V 137 matrices for the encoder self-attention layers as $Q = XW_O$, $K = XW_K$, and V =138 XW_V , where, $X = [\mathbf{y}_{k-w:k+P}; \mathbf{u}_{k-w:k+P}] \in \mathbb{R}^{b \times w+P \times l}$. The b and l denote the 139 batch size and the number of variables in the input matrix, respectively. Note that 140 the future values in sequence y in the X matrix $(y_{k+1:k+P})$ need to be masked with 141 the y value at the time step k (y_k) to avoid including the prediction output in the 142 input matrix (X) before training. 143

The prediction model, to be compatible with the MPC framework, requires 144 both future marked \mathbf{y}_k data points and $\mathbf{u}_{k+1:k+P}$ data points within the X matrix. 145 The $\mathbf{u}_{k+1:k+P}$ data points, in particular, are vital to ensure the degrees of freedom 146 for the MPC optimization solution. Conversely, while the MPC framework does 147 not inherently require future y points, they are essential for maintaining consistent 148 length between the y and u vectors within the X matrix. Including spurious y 149 values in the prediction horizon portion of the y vector can lead to anomalous 150 correlations between input and output. Nonetheless, this study demonstrates the 151 adeptness of Transformer networks in managing such data irregularities. 152

153 **3. Long-Short Term Memory Networks (LSTM)**

The LSTM network is a type of RNN along with a GRU (Gated Recurrent Unit) network. It is composed of a cell, an input gate, an output gate, and a forget gate to reinforce the long-term memory proposed in [21]. LSTMs are well suited for making time-series predictions with longer-term sequences and dependencies with this unique structure. However, the inherent sequential computation in the LSTM leads to an extended processing time that limits the online MPC application as a prediction model where a specific cycle time must be met. From a train-ability ¹⁶¹ point of view, the prediction accuracy with the capability of accommodating the

¹⁶² special input-out data structure for multistep-ahead prediction is compared to the

¹⁶³ proposed Transformer model. The LSTM block and equations for each gate are shown in Figure 3.



Figure 3: LSTM cell.

164

165 4. Transformer Model MPC Development

This section discusses the development of a Transformer model-based MPC system. It introduces the training data modification to obtain the multistep prediction model, layer configurations of the Transformer network, and an MPC framework that accommodates the Transformer prediction model.

170 4.1. Training Data Preparation

A simple single-input single-output (SISO) first-order plus dead-time (FOPDT) model is chosen for the process model for illustrative purposes. Equation 3 shows the continuous FOPDT equation between process input (*u*) and output (*y*), where K_p , τ_p , and θ_p represent the model parameters for process gain, time constant, and dead time, respectively. The parameter values used for the FOPDT model are also shown in Equation 3.

$$\tau_p \frac{dy(t)}{dt} = -y(t) + K_p u \left(t - \theta_p\right)$$
where, $K_p = 1, \tau_p = 2$, and $\theta_p = 0$
(3)

The training data is generated by simulating the FOPDT model with randomly 177 created input (u) sequences for 1,600 data points. The simulated data is split into 178 two sets, one for training and the other for validation. The validation set is used to 179 monitor and prevent over-fitting of the NN model. An additional approach to han-180 dle correlated multivariate control inputs is to build a Principal Component Anal-181 ysis (PCA) based model that reduces the input space or constrains the squared pre-182 diction error (SPE) to follow the PCA model during training. With unachievable 183 setpoints, the PCA-based SPE constraint helps to maintain achievable setpoints 184 [22, 23]. The data in this study does not have correlated inputs because they are 185 randomly generated. The training and validation data are reshaped to fit the time-186 series NN model structure. Although the Transformer and LSTM models use en-187 tirely different algorithms, the dimension of the input and output vector is the same 188 for the same sequence data. The three-dimensional vector, also called a tensor, can 189 be formed with (*Batchsize* \times *Length of one snapshot* \times *Number of variables*) and 190 the size of a tensor varies depending on the size of prediction steps such as one-191 step (os) and multistep (ms) ahead prediction. Figure 4 visually describes the 192 batch size or Number of snapshots (T), length of one snapshot (w+P), and num-193 ber of variables, respectively. The lengths of the look-back window (w = 5) and 194 the prediction horizon (P = 10) have been carefully selected to reflect the dynamic 195 characteristics of the process. The mathematical expressions of the data struc-196 ture are shown in Equation 4 and 5 for one-step- and multistep-ahead prediction 197 models, respectively. 198

$$Y_{os} = (y_{k+1}) \in \mathbb{R}^{b \times 1 \times l}$$

$$X_{os} = ((u_{k-w}, \cdots u_k), (y_{k-w}, \cdots y_k)) \in \mathbb{R}^{b \times w \times f}$$
where, $w \le k \le T - P$,
 $b = T - w - P$

$$Y_{ms} = (y_{k+1}, \cdots y_{k+P}) \in \mathbb{R}^{b \times P \times l}$$

$$X_{ms} = ((u_{k-w}, \cdots u_{k+P}), (y_{k-w}, \cdots y'_{k+1}, \cdots y'_{k+P}))$$
 $\in \mathbb{R}^{b \times (w+P) \times f}$
where, $w \le k \le T - P$,
 $b = T - w - P$,
 $y'_{k+i} = y_k \ (1 \le i \le P)$
(4)

¹⁹⁹ In Equations 4 and 5, *Y* and *X* represent the reshaped output and input of the NN ²⁰⁰ training, which are also called *labels* and *features*, respectively. On the other



Figure 4: Training data preparation with receding window snapshots

hand, u and y are the process system input and output of the FOPDT model. As 201 expressed in Equations 4 and 5, system input and output (u and y) are in the X 202 as a linear auto-regressive and exogenous part of the NN model input, respec-203 tively. The symbols b, T, w, and P denote the batch size, entire time-series length, 204 look-back window size, and prediction horizon length, respectively. Including all 205 snapshots k, l, and f denote the current time step, the number of labels, and the 206 number of features, respectively. A notable difference between one-step predic-207 tion and multistep prediction data structures is that the Y_{ms} extends the vector to 208 the k + P time step, including the full P-step-ahead values. In contrast, the Y_{os} in-200 cludes only the first value in the prediction horizon. A similar difference appears 210 in the input feature. The u and y vector in the X_{ms} extend to the k + P, while the 211 ones in the X_{os} stop at the k. Important customization needs to be made for the 212 y vector in the X_{ms} . The y values for future time steps are replaced with the y_k 213 shown as y'_{k+i} in the Equation 5. Prior state values are included in the data row 214 for training to maintain the vector shape consistent with the *u* vector in the X_{ms} . 215 Figure 5 visually compares the data structures between one-step- and multistep-216 ahead models for one snapshot sample in the batch. The modified y vector for 217 the multistep prediction model is illustrated in the dashed box with y_k . Figure 6 218 illustrates the training data modification for the multistep prediction models. 219

220 4.2. Transformer Model Training

Prepared training data sets for one-step-ahead prediction (Y_{os}, X_{os}) and multistepahead prediction (Y_{ms}, X_{ms}) are used for training the Transformer model. Every



Figure 5: Training data structure comparison between one-step-ahead model and multistep-ahead model

training result is compared to the result of the counterpart LSTM model. The 223 Transformer model consists of two encoders. The number of encoders and de-224 coders for each system can be determined by hyperparameter optimization. Two 225 encoders were used in this study. Hyperparameter optimization was not used 226 in this study, but could be the subject of future work. Each encoder includes a 227 multi-head attention layer, feedforward layer, dropout layer, and output layer. The 228 multi-head attention layer consists of 10 heads with softmax activation function 229 as described in Equation 2. The pre-processed inputs X_{os} and X_{ms} are used for 230 O, K, V vectors for the self-attention mechanism. The attention score vector is 231 concatenated with the input (X) to pose the residual to feed it to the following 232 feedforward layer. One feedforward layer consists of 100 neurons in the hidden 233 layer with a hyperbolic tangent (tanh) activation function, and the other output 234 feedforward layer with the same number of neurons with the number of variables 235 in feature X, which is 2 for this study. A dropout layer separates each feedforward 236 layer with a dropout factor of 20%. The attention score vector, which also repre-237 sents the output of the two encoder blocks, is fed to the final output feedforward 238 layer to map it to the size of the label Y: 1 for the one-step-ahead model and P for 230 the multistep-ahead model. 240



Figure 6: Training data modification for multistep-ahead prediction models

Meanwhile, the LSTM model consists of three LSTM layers and a dense layer, 241 each separated by a dropout layer with a dropout factor of 20%. Each LSTM layer 242 consists of 100 hidden states, and the current input and the 99 hidden states from 243 the previous time step are fed into the LSTM layers at any time step. The final 244 dense layer takes the 100 outputs generated by the third LSTM layer and maps 245 them to the outputs. Both networks are trained using an Adam optimizer, with loss 246 calculated as a mean-squared error (MSE). The summaries for multistep models 247 are shown in Table 1 and 2. The number of trainable parameters of weight and 248 bias is significantly less in the Transformer model than in LSTM. 249

Layer	Input Dimension	Output Dimension	Act. fn
MHA	$Batch \times 15 \times 2$	$Batch \times 15 \times 2$	SoftMax
Dense	$Batch \times 15 \times 2$	$Batch \times 15 \times 100$	tanh
Dense	$Batch \times 15 \times 100$	$Batch \times 15 \times 2$	Linear
MHA	$Batch \times 15 \times 2$	$Batch \times 15 \times 2$	SoftMax
Dense	$Batch \times 15 \times 2$	$Batch \times 15 \times 100$	tanh
Dense	$Batch \times 15 \times 100$	$Batch \times 15 \times 2$	Linear
Dense	$Batch \times 15 \times 2$	$Batch \times 10 \times 1$	Linear
Parameters		1,758	

Table 1: Summary of the Transformer model

Table 2: Summary of the LSTM model

Layer	Input Dimension	Output Dimension	Act. fn
LSTM	$Batch \times 15 \times 2$	$Batch \times 15 \times 100$	tanh
LSTM	$Batch \times 15 \times 100$	$Batch \times 15 \times 100$	tanh
LSTM	$Batch \times 15 \times 100$	$Batch \times 15 \times 100$	tanh
Dense	$Batch \times 15 \times 100$	$Batch \times 10 \times 1$	Linear
Parameters		203,010	

250 4.3. Transformer-Based Surrogate MPC

Two main parts of MPC are the prediction model and the optimization solver. In the NN-based surrogate MPC, NN models replace the classical transfer function model or physics-based model that generally uses differential equations. The optimization solver minimizes the MPC objective function by searching the best possible *u* sequences in the control horizon (*M*). The sum of squared error (*SSE*) objective function in the general MPC is shown in Equation 6.

$$\min_{\Delta U} \Phi = (\hat{Y} - Y_{ref})^T Q (\hat{Y} - Y_{ref}) + \Delta U^T R \Delta U$$
s.t. $0 = f(\dot{x}, x, y, p, d, u)$
 $0 = g(x, y, p, d, u)$
 $0 \le h(x, y, p, d, u)$
(6)

where \hat{Y} and Y_{ref} are the column vectors for model prediction values and the reference trajectory from the time step (k+1) to (k+P). ΔU is the column vector for the control moves for the future control horizon (*M*), from the time step (k+1)to (k+M). *Q* and *R* are the diagonal weighting matrices for multiple CVs (Controlled variables) and MVs (Manipulated variables). The detailed descriptions of the MPC equations can be found in [24].

The trained Transformer models are embedded into the MPC framework serving as prediction models providing the \hat{Y} in Equation 6. One-step-ahead prediction models for both LSTM and Transformer execute the models *P* times to get the full prediction for the defined prediction horizon. However, the multistepahead prediction models need only one execution for the same task, significantly improving computational efficiency.

5. Case Studies for Transformer Multistep-ahead Prediction Model MPC

This section discusses three MPC case studies with multistep Transformer 270 models. The case studies include two experiments in the simulation environments 271 and one with an actual temperature control test device, TCLab. The first case study 272 examines one batch of MPC optimization calculation results for a SISO FOPDT 273 process. The second and third case studies discuss continuous multi-input and 274 multi-output (MIMO) MPC tests for a fluidized-bed roaster in gold ore recovery 275 and with the TCLab device. The case studies contrast the improved computational 276 efficiency of the proposed model among the different types of networks (LSTM) 277 and prediction structures (one-step and multistep prediction). 278

279 5.1. Case I: First Order Dynamics Model (FOPDT)

Various types and structures of NN models are tested in the surrogate model MPC framework serving as prediction models, and the results are shown in Fig-

ure 7. Trained LSTM and Transformer models with one-step- and multistep-282 ahead prediction structures provide the model prediction values to the optimiza-283 tion solver. The gradient-based optimization solver iteratively searches the min-284 imum objective function value (SSE). A sequential least-squares programming 285 (SLSQP) solver is used in this study. The surrogate prediction models compute 286 multiple prediction results for every control interval. Figure 7 shows the optimally 287 predicted CV values based on the given setpoint of 1 and the computation time for 288 different models for one control interval. The plot shows the result of one con-289 trol interval in Figure 7 with the time steps in the prediction horizon (P = 10). 290 Compared to the one-step-ahead models, there is a significant improvement in 291 computation time for both LSTM and transformer with multistep-ahead predic-292 tion models. The multistep Transformer model shortens the computation time by 293 about 20 times compared to the one-step LSTM model (from 10.85s to 0.53s). 294 Computation time in RNN-based deep learning models is substantially decreased 295 by implementing two key modifications. Firstly, the sequence computation is re-296 placed with a self-attention mechanism, thereby eliminating the recurrence inher-297 ent in the one-step ahead LSTM model. Secondly, a multistep ahead prediction 298 structure is applied further enhancing the computational efficiency by calculating 299 the whole length of the prediction simultaneously. The self-attention mechanism 300 in the Transformer model accomplishes the same task with matrix operations re-301 ducing the internal sequences by w times less than the LSTM model. The other 302 significant portion that contributes to the computation time is the simultaneous 303 multistep-prediction structure. The multistep prediction structure simultaneously 304 computes the model prediction values for the entire prediction horizon (from k + 1305 to k + P), which reduces the number of model calls by P times less than the one-306 step prediction model. The overall number of sequential computations occurring 307 in one MPC control interval is shown in Table 3. 308

The number of iterations for MPC calculation (*n*) in Table 3 can be varied be-309 tween different models as the numerical solver converges differently with every 310 model. The traditional ODE (Ordinary Differential Equation) model-based MPC 311 is also tested along with the NN models to provide the baseline performance. For 312 a simple FOPDT model, the traditional MPC with a transfer function model out-313 performs NN models. As system complexity increases, similar NN models are 314 capable of capturing complex, higher-order relationships. In contrast, the effort to 315 build and the computation cost to execute physics-based models increases signif-316 icantly. The use of multistep Transformers to model complex, nonlinear systems 317 is a natural extension of this work and a topic of ongoing research. 318

The model accuracy is inferred from Figure 7. First, the u values are supposed

Model	Sequential	Multistep	MPC	Number of
Туре	Computation	Prediction	Iteration	Executions
ODE (Sequential)	1	Р	п	$P \times n$
LSTM-OS	W	Р	n	$w \times P \times n$
LSTM-MS	w + P	1	n	$(w+P) \times n$
Transformer-OS	1	Р	n	$P \times n$
Transformer-MS	1	1	n	n

Table 3: Number of sequential computations occurring in one MPC calculation, where, w, P, and n represent a look-back window, prediction horizon, function evaluation of MPC optimization

to be at a steady-state '1.0' like the ODE model result because the process gain 320 (K_p) in the FOPDT model is defined as '1.0'. The steady-state u value must be 321 the same as the y value with '1.0' as the output setpoint. However, the u results 322 of the multistep-ahead LSTM model and the one-step-ahead Transformer model 323 converge to slightly different values than 1.0, which means those models have 324 some model mismatch with the process model. In addition, the LSTM multistep 325 model shows a significant model mismatch on both y and u, while the multistep 326 Transformer model follows the ODE model result. The different outcomes can be 327 explained by two distinguished methods to learn the temporal dependency. The 328 suggested training data structure for the multistep-ahead prediction models in-329 cludes both the past measured data (k - w : k) and future prediction horizon data 330 (k+1: k+P), especially for the feature (X) set. As discussed in Section 4.1, this 331 is necessary not only for including the entire prediction horizon data to learn si-332 multaneously but also to pass the control horizon part of the feature $(u_{k+1}: u_{k+M})$ 333 to the SLSQP solver as decision variables to be solved. However, while the pre-334 diction horizon is part of the data in the feature set (X), the prediction horizon 335 part of y data $(y_{k+1} : y_{k+P})$ is replaced with the currently measured data (y_k) . This 336 modification avoids including the labeled data in the feature set before training 337 while maintaining a consistent data dimension. However, irrelevant data is in-338 cluded due to this modification, which introduces irregular temporal dependency 339 for time-series data training. The way that the Transformer architecture learns 340 the temporal dependency is completely different than RNNs. The self-attention 341 mechanism employs a conditional probability, selectively including the useful re-342 lationship and excluding the irrelevant relationship between input and output data. 343 Therefore, the Transformer is more effective in learning the irregular temporal 344

dependency than LSTM that structurally includes the temporal dependency in a
 sequential manner.

347 5.2. Case II: Fluidized-Bed Roaster Model

The roasting process is a part of gold production that recovers the gold from 348 the ore. The substances that negatively affect the carbon-in-leach (CIL) process 349 are removed by the oxidation reaction in the roasting process. The particular type 350 of gold ore, called refractory gold ore, has several characteristics that interfere 351 with satisfactory gold recovery in the traditional CIL. The iron sulfide in the re-352 fractory gold ore confines the gold in the minerals that prevent contact with the 353 cyanide solution. The carbonaceous minerals in the ore tend to re-absorb the re-354 covered gold component from the cyanide solution. Thus, the refractory gold ore 355 needs to be pretreated in the roaster, converting the iron sulfide mineral and car-356 bonaceous mineral to oxides. The roaster consists of two stages of fluidized bed 357 reactors (FBRs) that oxidize sulfide compounds and organic carbon in the ore. 358 Before the roasting process, the ore is crushed into fine particles with an average 359 size of $75\mu m$ in diameter by passing them through a grinding process. The fine 360 particles are fed into the roaster from the top and are fluidized by the upward force 361 generated by the oxygen flow from the bottom of the roaster. Three critical factors 362 that affect the oxidation reaction are ore particle size, oxygen content in the com-363 bustion gas, and bed temperature. Compared with traditional air-roasting, oxygen-364 roasting yields better gold recovery by improving reaction efficiency. Compared 365 to oxygen roasting, air roasting requires a longer retention time at a higher tem-366 perature to maintain a high conversion of sulfidic metals and carbonaceous matter. 367 These conditions easily over-roast and soften the ore particles (glassy flux), block-368 ing the pores of the material from contacting the combustion gas. Over-roasting 369 also affects the subsequent CIL process by encapsulating the gold in the glassy 370 fluxed ore [25, 26, 27, 28]. The exothermic reaction generates most of the thermal 371 energy required for the oxidation reaction. Additional heat is provided by feeding 372 sulfur prills as fuel. The major reactions [27] are shown below and a diagram of 373 the roasting process is shown in Figure 8. 374

Combustion of organic carbon:

$$C(s) + O_2(g) \rightarrow CO_2(g)$$



(a) Setpoint tracking performance of CV and computation time



(b) Control movement of MV corresponding to given setpoint

Figure 7: Results of MPC calculation compared to the various types and structures of NN models. Setpoint tracking performance of CV and computation time (Upper), and control movement of MV corresponding to given setpoint (Lower)



Figure 8: Schematic drawing of the roaster process

Combustion of iron sulfides (pyrite and pyrrhotite):

$$7 \operatorname{FeS}_2(s) + 6 \operatorname{O}_2(g) \to \operatorname{Fe}_7 \operatorname{S}_8(s) + 6 \operatorname{SO}_2(g)$$

 $4 \operatorname{Fe}_7 \operatorname{S}_8(s) + 53 \operatorname{O}_2(g) \to 14 \operatorname{Fe}_2 \operatorname{O}_3(s) + 32 \operatorname{SO}_2(g)$

Combustion of sulfur prills (fuel):

$$S(s) + O_2(g) \rightarrow SO_2(g)$$

Retention of Sulfur dioxide:

$$2\operatorname{CaCO}_3(s) + 2\operatorname{SO}_2(g) + \operatorname{O}_2(g) \rightarrow 2\operatorname{CaSO}_4(s) + 2\operatorname{CO}_2(g)$$

The actual operation data is usually insufficient to train the NN surrogate mod-375 els because the actual operations are often maintained in a narrow range that does 376 not provide the information for the broader operation range where the optimal 377 operation point may exist. Thus, a physics-based model is developed for process 378 simulation experiments, including data generation for training NN models and 379 emulating the actual process as a digital twin model during the controller perfor-380 mance test. The physics-based model for this research consists of heat and mass 381 balance equations and reaction kinetics. The reaction kinetics uses a shrinking 382 core model (SCM), whose reaction rate is dominated by pore diffusion. 383

³⁸⁴ 5.2.1. Physics-based Model of a Roaster Process

A physics-based model for the roaster is developed to generate the Trans-385 former model training and validate the Transformer model-based MPC perfor-386 mance. The SCM is used for reaction kinetics of ore particle oxidation reaction 387 shown in Equation 7. The rate equation of the SCM consists of three parts: Diffu-388 sion through gas film controls, diffusion through ash layer control, and chemical 389 reaction controls. The diffusion through gas film control is negligible for engi-390 neering applications because it has an effect only for a short period of the initial 391 stage. Once the ash layer starts forming on the particle surface, the ash layer dif-392 fusion step controls the overall reaction rate [29, 30]. The ash layer diffusion term 393 requires a calculation of the effective diffusivity, \mathcal{D}_e . The Chapman-Enskog equa-394 tion is used for calculating the effective diffusivity (\mathcal{D}_e) shown in Equations 8 -395 10. 396

$$-r_{i} = \frac{A_{p}C_{i}}{\frac{1}{k_{g}} + \frac{R_{p}}{6\mathscr{D}_{e}} + \frac{1}{k''}}$$

$$= \frac{A_{p}C_{i}}{\frac{R_{p}}{6\mathscr{D}_{e}} + \frac{1}{k''}}$$
(7)

$$\mathscr{D}_{AM} = 0.0018583 \sqrt{T^3 \left(\frac{1}{M_{O_2}} + \frac{1}{M_{O_2}}\right)} \frac{1}{p \sigma_{AB}^2 \Omega_{\mathscr{D},AB}}$$
(8)

$$\Omega_{\mathscr{D},AB} = \frac{1.06036}{T^{*0.15610}} + \frac{0.19300}{\exp(0.47635T^*)} + \frac{1.03587}{\exp(1.52996T^*)} + \frac{1.76474}{\exp(3.89411T^*)}$$

where, $T^* = \kappa T/\varepsilon$ (9)

$$\mathscr{D}_e = \frac{\mathscr{D}_{AM}\varepsilon_p}{\tau}, \ where \ \tau = \varepsilon_p^{-0.41}$$
 (10)

$$k'' = A_p T \exp\left(-E_r/RT\right) \tag{11}$$

$$\frac{dN_i}{dt} = F_{i,in} - F_{i,out} + r_i \tag{12}$$

$$NC_p \frac{dT_i}{dt} = F_{i,0}H_{in} - F_{i,out}H_{out} + \Delta H_{reaction}$$
(13)

Table 4: Lumped Parameters from Physics-based Roaster Model

Quantity	Value
Average surface area of ore particle (A_p)	$0.0176 m^2$
Average radius of ore particle (R_p)	37.5 µm
Atmospheric pressure (p)	1 <i>atm</i>
Molecular weight of Oxygen (M_{O_2})	32 g/mol
Lennard-Jones potential parameter for Oxygen molecules (σ)	3.433 Å
Lennard-Jones potential parameter for Oxygen molecules (ε/K)	113 K
Average porosity of ore particles (ε_p)	0.5
Effective diffusivity (\mathcal{D}_e)	
Diffusivity for gases (\mathcal{D}_{AM})	
Mass transfer coefficient between air and particle (k_g)	
Rate constant for the surface reaction (k'')	
Reactor temperature (T)	
Tortuosity (τ)	
Collision Integral for diffusivity (Ω)	

³⁹⁷ 5.2.2. Train Transformer Model for Roaster MPC

Training data is generated by simulating the physics-based roaster model. The 398 roaster simulation model has fourteen main operation variables, including six in-399 put and eight output variables. The first set of input variables is related to the 400 quantity of reactor feed, such as the quantity of ore feed, the mass flow rate of 401 sulfur prill, and the volume flow rate of oxygen. These are the primary operation 402 variables of roaster operation and are the candidates for an advanced process con-403 trol system manipulated variables (MV). The second set of input variables is the 404 quality of the ore feed, such as the contents of a specific component. Compared 405 to the first set of inputs, the second set is not adjustable for the process operation; 406 thus, it is considered an advanced process control system disturbance variable 407 (DV). The roaster model in this research includes organic carbon, iron sulfide, 408 and carbonate content in the feed ore. The random input signals for the six reactor 409 input variables are generated and simulated in the reactor model, which creates 410 the associated output variable responses. The random input signals and output 411 variable responses are then used for training the Transformer model. The reactor 412 output variables in this case study include the off-gas and calcine specifications 413 and the two reactor temperatures. The input and output variables of the roaster 414 process are shown in Table 5 and the roaster operation range for the random input 415 signals refer to [27]. Fifteen thousand sample points are generated for training the 416 Transformer model, and a short segment of the input and output data is shown in 417 Figures 9 and 10. Figure 10 also displays the Transformer model fitting. The blue 418 dashed line in the figure represents the sampled data from the roaster physics-419 based simulation model, and the red line shows the prediction results with the 420 Transformer model. The model fitting results validate that the Transformer model 421 can provide accurate model predictions for control. 422

The network configuration for the Transformer model is shown in Table 6. It 423 consists of Multi-head attention (MHA) layers followed by a set of Feedforward 424 layers, as depicted in Figure 2 in the previous section 2. The tensor dimension for 425 each layer represents: $(Batch) \times (Length of one snapshot) \times (Number of variables)$, 426 where the *Batch* size for the roaster training is the same as the length of the time 427 series. Each snapshot includes the past data with a length of look-back window 428 size, of six, plus the prediction horizon length, of ten, which forms the second 429 dimension of MHA input tensor, sixteen. The output dimension of the last feed-430 forward layer (*Dense* layer) is reduced to the size of the multistep prediction result 431 with the length of the prediction horizon, ten, and the number of system output 432 variables, eight: $(Batch \times 10 \times 8)$. 433



Figure 9: Roaster input data for Transformer model training



Figure 10: Transformer model training result for Roaster



Figure 11: Transformer model validation result for Roaster

Variables in Roaster model			
Input variables	Output variables		
	O ₂ in Offgas (wt%)		
Ore feed amount (Amp)	CO_2 in Offgas (wt%)		
Sulfur prill (T/H)	SO ₂ in Offgas (wt%)		
Oxygen (SCFM)	TCM in Calcine (wt%)		
Carbon content in Ore (wt%)	<i>FeS</i> ² in Calcine (wt%)		
Iron Sulfide content in Ore (wt%)	<i>CaCO</i> ₃ in Calcine (wt%)		
Carbonate content in Ore (wt%)	Stage1 Temperature (F)		
	Stage2 Temperature (F)		

Table 5: Operation variables of roaster simulation model

Table 6: Summary of the Transformer model for roaster data

Layer	Input Dimension	Output Dimension	Act. fn
MHA	$Batch \times 16 \times 14$	$Batch \times 16 \times 14$	SoftMax
Dense	$Batch \times 16 \times 14$	$Batch \times 16 \times 100$	tanh
Dense	$Batch \times 16 \times 100$	$Batch \times 16 \times 14$	Linear
MHA	$Batch \times 16 \times 14$	$Batch \times 16 \times 14$	SoftMax
Dense	$Batch \times 16 \times 14$	$Batch \times 16 \times 100$	tanh
Dense	$Batch \times 16 \times 100$	$Batch \times 16 \times 14$	Linear
Dense	$Batch \times 16 \times 14$	$Batch \times 10 \times 8$	Linear
Parameter		26,216	

434 5.2.3. Control Results for Roaster

An MPC controller is designed as a case study to test the prediction accu-435 racy and solution time of the Transformer based model in the MPC system. This 436 case study selects two controlled variables (CVs) and two manipulated variables 437 (MVs). The two reactor temperatures at each stage of the roaster are chosen for 438 two CVs, and ore feed rate and sulfur prill inlet flows are chosen for the two MVs. 439 The reactor temperature is one of the most critical operation variables directly af-440 fected by the reaction status. Thus, it is commonly selected as a major CV in many 441 advanced process control applications inferring the conversions. The ore particles 442 are transported by the bucket elevator to the 1st stage roaster, and the feed rate is 443 adjusted by changing the speed of the bucket elevator. The sulfur prill feed rate is 444 assigned to the second MV that maintains the reaction temperature requirement to 445

ensure high sulfide sulfur and organic carbon conversion. The other process input
variables, including the ore feed composition and oxygen inlet flow, are retained
at a constant value during the experiment. The MPC settings are shown in Table
7.

Darameters	Setting
	Setting
Prediction Horizon (P)	10 min
Control Horizon (<i>M</i>)	6 min
CV Weight for SSE (q_{T1})	$1x10^{3}$
CV Weight for SSE (q_{T2})	$1x10^{2}$
MV Weight for rate of change $(r_{Ore feed})$	$1x10^{3}$
MV Weight for rate of change $(r_{Sulfur prill})$	1
MV Max move (Ore feed)	1 Amp
MV Max move (Sulfur prill)	0.2 T/H
MV Bounds for SLSQP (Ore feed)	80 - 100 <i>Amp</i>
MV Bounds for SLSQP (Sulfur prill)	9 - 13 <i>T /H</i>

Table 7: Settings of Roaster MPC

A simulated MPC control result for the roaster process is shown in Figure 12. 450 In this test, two temperature setpoints are changed to validate the setpoint tracking 451 performance of the Transformer model-based MPC. At 20 minutes, MPC is ac-452 tivated and starts controlling the temperatures. Ore and sulfur feed rates as MVs 453 move to meet the temperature setpoints based on the MPC control calculation. To 454 increase the temperatures at a given condition, MPC decreases ore feed and in-455 creases the sulfur prill feed because the sulfur prill has a higher reaction rate than 456 ore particle as it does not go through the ash layer diffusion. This demonstrates 457 that the Transformer controller provides proper prediction by taking into account 458 the reaction kinetics of gold ore roasting. Although this case study only includes 459 a subset of operation variables, the complete set of variables can be included in 460 future research studies not only for the MPC application but also for the real-time 461 optimization application that deals with an economic objective and steady-state 462 targets. 463

464 5.3. Case III: Temperature Control Lab Device

In this section, an experiment is conducted using a temperature control lab (TCLab). The device consists of two heater inputs and two temperature output measurements. This device has been used for control education and research



Figure 12: MPC result for Roaster temperature control

⁴⁶⁸ purposes, demonstrating control theory and machine learning-based automation ⁴⁶⁹ methods [31].

In the experiment, the new Transformer multistep-ahead prediction model is 470 developed for the TCLab device and embedded into the MPC algorithm. The 471 newly proposed Transformer multistep-ahead prediction model is compared with 472 three other options, including the Transformer one-step-ahead prediction model 473 and one-step and multistep-ahead prediction models with the LSTM network. In 474 the TCLab case study, all four different model types are tested for offline model 475 development to show model fitness. Then, two of the four models are chosen 476 to validate the online control performance: the multistep Transformer model and 477 the one-step LSTM model. The two models are compared for the computational 478 efficiency and setpoint tracking performance in the online control comparison. 479

The training data is generated with the TCLab device with random step test 480 input signals for the two heaters. The step-test random signals are generated with 481 a combination of two individual random number generators. The first regulates 482 the amplitude of step changes, while the second generates the number of intervals 483 between the step changes. The amplitude range for the first random number gen-484 erator is zero to one hundred percent of the maximum heater value of the TCLab. 485 The step interval range for the second random number generator is in the range of 486 three hundred to six hundred seconds to account for the dynamic behavior of the 487 TCLab. This step-change interval distribution ensures the collection of various 488 input-output relationships with variable frequency responses. 489

Data is collected every second for twenty-five thousand seconds and down-490 sampled every 30 seconds balancing the control calculation speed and control 491 performance in the MPC application. The first training data set includes the first 492 3,000 seconds of the data to represent a case of insufficient training data. The 493 other training session uses the full 25,000 seconds of data. The validation results 494 with MAE values for each model are shown in Table 8. In addition, our case stud-495 ies have certain limitations when confirming the LSTM model performance with 496 irregular temporal dependencies. While this has been repeatedly demonstrated 497 in the field of natural language processing, the model prediction accuracy shown 498 in Table 8 does not fully encapsulate the MS LSTM viability as an MPC pre-499 diction model. Thus, a more extensive study with a primary focus on prediction 500 performance could be an interesting topic for future research. Figure 13 depicts 501 the model training results of using either a small data set or large data to train 502 one-step LSTM models and multistep Transformer models on TCLab data. Fig-503 ure 13a, shows the one-step model training result using the small three thousand 504 seconds of data for training. The model predicts the general shape of the mea-505

sured data but is not able to accurately predict the measured values. Figure 13b 506 shows the training result of the one-step LSTM model trained on the large twenty-507 five thousand seconds of data. The model accurately predicts both the shape of 508 the measured data and the values of the data over the training range. Figure 13c 509 shows the training result for the Transformer one-step model using three thousand 510 seconds of data. Compared with the LSTM one-step model trained in Figure 13a 511 on the smaller data range, overall accuracy is greatly improved. In addition, when 512 Figure 13c is compared with the LSTM model trained on the large set of data 513 depicted in Figure 13b, the Transformer model, despite having less training data, 514 is more accurate in certain intervals than the LSTM. Figure 13d shows the results 515 for using the twenty-five thousand-second data set on the Transformer one-step 516 model. Of the four models in Figure 13, it is best able to predict the measured 517 data across the entire range. 518



(a) LSTM (one-step model): using 3,000 seconds of data



(c) Transformer (one-step model): using 3,000 seconds of data)



(b) LSTM (one-step model): using 25,000 seconds of data)



(d) Transformer (one-step model): using 25,000 seconds of data)

Figure 13: One-step Models validation for LSTM and Transformer: in comparison with using 3,000 (a, c) and 25,000 seconds data (b, d)

⁵¹⁹ Figure 14 depicts the results of using either a small data set or large data to



(a) LSTM (Multistep model): using 3,000 seconds of data



80 LSTM Measured 티 60 \sim 40 4000 6000 8000 10000 14000 12000 LSTM 60 Measured 입 20 V 40 30 14000 12000 4000 6000 8000 10000

(b) LSTM (Multistep model): using 25,000 seconds of data)



(c) Transformer (Multistep model): using 3,000 seconds of data)

(d) Transformer (Multistep model): using 25,000 seconds of data)

Figure 14: Multi-step Models validation for LSTM and Transformer: in comparison with using 3,000 (a, c) and 25,000 seconds data (b, d)

			M	AE
			Quantity of	training data
Prediction Type	Network Type	Variables	3,000 sec	24,000 sec
One-step	LSTM	T_1 (° C)	4.03	3.82
		T_2 (° C)	3.88	2.64
	Transformer	T_1 (° C)	2.16	3.32
		T_2 (° C)	5.26	5.63
Multistep	LSTM	T_1 (° C)	3.59	2.59
		T_2 (° C)	2.70	1.57
	Transformer	T_1 (° C)	2.89	2.61
		T_2 (° C)	2.87	1.47

Table 8: Training result comparison

train multi-step LSTM models and Transformer models on TCLab data. Figure 520 14a shows the results of using the smaller three thousand second data to train a 521 multistep LSTM model. Figure 14b depicts the results of using the larger twenty-522 five thousand data set to train the multistep LSTM model. The model accurately 523 predicts the measured data. Figure 14c show the results of training a multistep 524 Transformer model on the small three thousand second data. It is more accurate 525 than the multistep LSTM trained on 3,000 seconds of data shown in Figure 14a, 526 and in certain ranges, it is more accurate than the multistep LSTM model trained 527 on twenty-five thousands seconds of data shown in Figure 14b. Figure 14d depicts 528 the result of using the twenty-five thousand seconds data set to train the multistep 520 Transformer model. This model is the most accurate of the four models, as ex-530 pected. 531

The LSTM one-step prediction model and the Transformer multistep predic-532 tion model are chosen to validate the online control performance, as shown in Fig-533 ure 15. As observed in the FOPDT model experiment in Section 5.1, a significant 534 difference in MPC calculation time is observed between LSTM one-step model 535 and the Transformer multistep model. The solution time at each control interval 536 is also shown in Figure 15a and 15b, along with the TCLab system variables. 537 The average solution time of the Transformer multistep model is 1.64 seconds. 538 In contrast, the LSTM one-step model takes 28.5 seconds on average, which is 539 more than fifteen times slower than the Transformer multistep model. The MPC 540 solution time affects the control performance in many different ways. The slower 541 solution time must be considered to select the control interval during the MPC 542 design. In this test, for example, the control interval for the Transformer multistep 543 model can be set to 2 seconds to account for overhead communication time and 544 variable solver iterations. However, the LSTM one-step model control interval is 545 set to 30 seconds as it needs longer than 28.5 seconds of average MPC solution 546 time. This slower MPC solution time also affects the overall control performance. 547 The mean absolute error (MAE) value between setpoint and temperatures mea-548 sures the setpoint tracking performance of each type of controller and is observed 549 visually in Figure 15a and 15b. Note that the total simulation time to fulfill the 550 identical setpoint sequences for both controllers is different: sixty-eight hundred 551 seconds for LSTM one-step model and thirty-six hundred seconds for the Trans-552 former multistep model. The surrogate MPC efficiency affects setpoint tracking 553 performance or disturbance compensating performance. For example, the rise 554 time, the amount of time that the controlled variable (CV) reaches the setpoint for 555 the first time since the latest setpoint has been changed, is 100 seconds, while the 556 rise time of the LSTM takes more than 500 seconds on average. The weights for 557

the MPC objective function for both controllers are set to the same settings. Table9 summarizes the control performance metrics.

Model Type		LSTM One-step	Transformer Multistep
Simulation time		6885s	3600s
Control interval		30s	10s
Number of MPC executions		115	345
	Min	13.61s	0.79s
MPC solution time	Max	51.07s	3.60s
	Avg	28.54s	1.64s
MAE (Mean absolute error)	T1 (° <i>C</i>)	10.97	8.26
	T2 (° <i>C</i>)	11.03	8.01

Table 9: MPC performance comparison between LSTM and Transformer models



(b) TCLab MPC result: Transformer (Multistep model)

Figure 15: TCLab MPC control result comparison between one-step LSTM model and Multistep Transformer model

560 6. Conclusions

A novel Transformer with a multistep prediction model structure is proposed 561 for MPC. The new multistep prediction model structure is tested with the Trans-562 former NN architecture. The Transformer model demonstrates improved com-563 putation time with the multistep prediction structure compared to the one-step 564 prediction structures with the LSTM model. The parallel data processing of the 565 Transformer architecture eliminates the sequential computation nature of RNN-566 based models, including LSTM. Additionally, the simultaneous multistep predic-567 tion structure removes the recursive procedure of the one-step prediction model 568 that updates the prediction horizon value one step at a time. As a result, the 569 proposed multistep Transformer model calculates the entire length of the model 570 prediction with a single execution of the forward propagation. The increased com-571 putational speed of a control system, such as with MPC, substantially improves 572 the control performance by enabling a shorter control interval in real-time control 573 practice. 574

In addition, the self-attention mechanism in the Transformer model learns 575 the irregular temporal dependencies more effectively than the sequential learn-576 ing mechanism in the LSTM. The irregular temporal dependencies introduced 577 in the multistep prediction structure are properly handled with the Transformer 578 model compared to the LSTM, which is briefly shown in the first case study result 579 (Figure 7). The multistep Transformer model shows more dependable accuracy 580 than the multistep LSTM model, even with a smaller model size, containing fewer 581 trainable parameters (Table 1 and 2). 582

This study does not carry out a thorough evaluation of the model quality across a variety of problem complexities or examine how different types of models, including alternatives like MS-LSTM and OS-Transformer, would perform under these conditions. Such an investigation could serve as an intriguing focus for future research endeavors.

588 **References**

- [1] J. H. Lee, J. Shin, M. J. Realff, Machine learning: Overview of the recent progresses and implications for the process systems engineering field, Com-
- ⁵⁹¹ puters and Chemical Engineering 114 (2018) 111–121. doi:10.1016/j. ⁵⁹² compchemeng.2017.10.008.
- ⁵⁹³ URL http://dx.doi.org/10.1016/j.compchemeng.2017.10.008

- [2] S. J. Qin, L. H. Chiang, Advances and opportunities in machine learning
 for process data analytics, Computers and Chemical Engineering 126 (2019)
 465–473. doi:10.1016/j.compchemeng.2019.04.003.
- [3] V. Venkatasubramanian, The promise of artificial intelligence in chemical
 engineering: Is it here, finally?, AIChE Journal 65 (2) (2019) 466–478. doi:
 10.1002/aic.16489.
- [4] R. B. Gopaluni, A. Tulsyan, B. Chachuat, B. Huang, J. M. Lee, F. Amjad, S. K. Damarla, J. Woo Kim, N. P. Lawrence, Modern machine learning tools for monitoring and control of industrial processes: A survey, IFAC-PapersOnLine 53 (2) (2020) 218–229. doi:10.1016/j.ifacol.2020.12.
 126.
- ⁶⁰⁵ URL https://doi.org/10.1016/j.ifacol.2020.12.126
- [5] S. W. Wang, D. L. Yu, J. B. Gomm, G. F. Page, S. S. Douglas, Adaptive neural network model based predictive control of an internal combustion engine with a new optimization algorithm, Proceedings of the Institution of Mechanical Engineers, Part D: Journal of Automobile Engineering (2006).
 doi:10.1243/095440706X72754.
- [6] M. Lee, S. Park, A new scheme combining neural feedforward control with
 model-predictive control, AIChE Journal 38 (2) (1992) 193–200. doi:10.
 1002/aic.690380204.
- [7] R. K. Al Seyab, Y. Cao, Differential recurrent neural network based pre dictive control, Computers and Chemical Engineering 32 (7) (2008) 1533–
 1545. doi:10.1016/j.compchemeng.2007.07.007.
- [8] Y. Wang, K. Velswamy, B. Huang, A long-short term memory recurrent
 neural network based reinforcement learning controller for office heating
 ventilation and air conditioning systems, Processes 5 (3) (9 2017). doi:
 10.3390/pr5030046.
- [9] J. Drgoňa, A. R. Tuor, V. Chandan, D. L. Vrabie, Physics-constrained deep
 learning of multi-zone building thermal dynamics, Energy and Buildings 243
 (2021). doi:10.1016/j.enbuild.2021.110992.
- ⁶²⁴ [10] W. C. Wong, J. Li, X. Wang, E. Chee, J. Li, X. Wang, Recurrent neural network-based model predictive control for continuous pharmaceutical man-

- ufacturing, Mathematics 6 (11) (7 2018). doi:10.3390/math6110242.
 URL http://arxiv.org/abs/1807.09556
- [11] Z. Wu, A. Tran, Y. M. Ren, C. S. Barnes, S. Chen, P. D. Christofides, Model
 predictive control of phthalic anhydride synthesis in a fixed-bed catalytic
 reactor via machine learning modeling, Chemical Engineering Research and
 Design 145 (2019) 173–183. doi:10.1016/j.cherd.2019.02.016.
- [12] J. Park, R. A. Martin, J. D. Kelly, J. D. Hedengren, Benchmark temperature microcontroller for process dynamics and control, Computers & Chemical Engineering 135 (2020) 106736.
 doi:10.1016/j.compchemeng.2020.106736.
- 636 URL https://linkinghub.elsevier.com/retrieve/pii/ 637 S0098135419310129
- [13] D. Machalek, J. Tuttle, K. Andersson, K. M. Powell, Dynamic energy
 system modeling using hybrid physics-based and machine learning en coder-decoder models, Energy and AI 9 (2022) 100172. doi:10.1016/
 j.egyai.2022.100172.
- [14] A. Mesbah, K. P. Wabersich, A. P. Schoellig, M. N. Zeilinger, S. Lucia,
 T. A. Badgwell, J. A. Paulson, Fusion of Machine Learning and MPC under
 Uncertainty: What Advances Are on the Horizon?, in: Proc. of the American
 Control Conference (ACC), 2022, pp. 342–357.
- [15] B. K. M. Powell, D. Machalek, T. Quah, Real-time optimization using re inforcement learning, Computers and Chemical Engineering 143 (2020)
 107077. doi:10.1016/j.compchemeng.2020.107077.
- ⁶⁴⁹ URL https://doi.org/10.1016/j.compchemeng.2020.107077
- [16] D. Machalek, T. Quah, K. M. Powell, A Novel Implicit Hybrid Machine
 Learning Model and Its Application for Reinforcement Learning, Computers & Chemical Engineering 155 (2021) 107496. doi:10.1016/j.
 compchemeng.2021.107496.
- ⁶⁵⁴ URL https://doi.org/10.1016/j.compchemeng.2021.107496
- [17] H. Yoo, H. E. Byun, D. Han, J. H. Lee, Annual Reviews in Control Rein forcement learning for batch process control : Review and perspectives, Annual Reviews in Control 52 (July) (2021). doi:10.1016/j.arcontrol.
 - 37

2021.10.006.
 URL https://doi.org/10.1016/j.arcontrol.2021.10.006

- [18] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez,
 L. Kaiser, I. Polosukhin, Attention is all you need, Advances in Neural In formation Processing Systems 2017-Dec (Nips) (2017) 5999–6009.
- [19] D. Bahdanau, K. H. Cho, Y. Bengio, Neural machine translation by jointly
 learning to align and translate, 3rd International Conference on Learning
 Representations, ICLR 2015 Conference Track Proceedings (2015) 1–15.
- [20] K. Cho, M. B. Van, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk,
 Y. Bengio, Learning Phrase Representations using RNN Encoder–Decoder
 for Statistical Machine Translation, Journal of Clinical Microbiology 28 (4)
 (2014) 828–829. doi:10.1128/jcm.28.4.828-829.1990.
- [21] S. Hochreiter, Long Short-Term Memory, Neural computation 1780 (1997)
 1735–1780.
- [22] H. Hassanpour, B. Corbett, P. Mhaskar, Integrating dynamic neural network models with principal component analysis for adaptive model predictive control, Chemical Engineering Research and Design 161 (2020) 26–37.
 doi:https://doi.org/10.1016/j.cherd.2020.03.031.
- [23] H. Hassanpour, B. Corbett, P. Mhaskar, Artificial neural network-based
 model predictive control using correlated data, Industrial & Engineering
 Chemistry Research 61 (8) (2022) 3075–3090.
- [24] J. Park, C. Price, D. Pixton, M. Aghito, R. Nybø, K. Bjørkevoll, J. D. Heden gren, Model predictive control and estimation of managed pressure drilling
 using a real-time high fidelity flow model, ISA Transactions 105 (2020) 256–
 268. doi:10.1016/j.isatra.2020.05.035.
- ⁶⁸³ URL https://doi.org/10.1016/j.isatra.2020.05.035
- [25] C. W. Jr, Oxygen Engiched Atmosphere Roasting, Ph.D. thesis, Montana
 School of Mines (1948).
- ⁶⁸⁶ [26] J. C. Smith, T. H. McCord, G. R. O'Neil, Treating refractory gold ores via ⁶⁸⁷ oxygen-enriched roasting (1990).

- [27] K. Thomas, A. Cole, Roasting Developments Especially Oxygenated
 Roasting, Ken Thomas & Murray Pearson, 2016. doi:10.1016/
 b978-0-444-63658-4.00023-2.
- ⁶⁹¹ URL http://dx.doi.org/10.1016/B978-0-444-63658-4.00023-2
- [28] H. Qin, X. Guo, Q. Tian, D. Yu, L. Zhang, Recovery of gold from sulfide re fractory gold ore: Oxidation roasting pretreatment and gold extraction, Min erals Engineering 164 (January) (2021). doi:10.1016/j.mineng.2021.
 106822.
- [29] D. Kunii, O. Levenspiel, Fluidization engineering, Butterworth-Heinemann,
 1991.
- [30] H. Ahn, S. Choi, A comparison of the shrinking core model and the grain
 model for the iron ore pellet indurator simulation, Computers and Chemical
 Engineering 97 (2017) 13–26. doi:10.1016/j.compchemeng.2016.11.
 005.
- [31] J. Park, R. A. Martin, J. D. Kelly, J. D. Hedengren, Benchmark temperature microcontroller for process dynamics and control, Computers and Chemical Engineering 135 (4 2020). doi:10.1016/j.compchemeng.2020.106736.