# Multi-task Optimization Based Co-training for Electricity Consumption Prediction

Hui Song
*School of Engineering*
*RMIT University*
Melbourne, Australia
hui.song@rmit.edu.au

A. K. Qin
*Department of Computing Technologies*
*Swinburne University of Technology*
Hawthorn, Australia
kqin@swin.edu.au

Chenggang Yan
*School of Automation*
*Hangzhou Dianzi University*
Hangzhou, China
cgyan@hdu.edu.au

*Abstract*—Real-world electricity consumption prediction may involve different tasks, e.g., prediction for different time steps ahead or different geo-locations. These tasks are often solved independently without utilizing some common problem-solving knowledge that could be extracted and shared among these tasks to augment the performance of solving each task. In this work, we propose a multi-task optimization (MTO) based co-training (MTO-CT) framework, where the models for solving different tasks are co-trained via an MTO paradigm in which solving each task may benefit from the knowledge gained from when solving some other tasks to help its solving process. MTO-CT leverages long short-term memory (LSTM) based model as the predictor where the knowledge is represented via connection weights and biases. In MTO-CT, an inter-task knowledge transfer module is designed to transfer knowledge between different tasks, where the most helpful source tasks are selected by using the probability matching and stochastic universal selection, and evolutionary operations like mutation and crossover are performed for reusing the knowledge from selected source tasks in a target task. We use electricity consumption data from five states in Australia to design two sets of tasks at different scales: a) one-step ahead prediction for each state (five tasks) and b) 6-step, 12-step, 18-step, and 24-step ahead prediction for each state (20 tasks). The performance of MTO-CT is evaluated on solving each of these two sets of tasks in comparison to solving each task in the set independently without knowledge sharing under the same settings, which demonstrates the superiority of MTO-CT in terms of prediction accuracy.

*Index Terms*—Multi-task optimization, inter-task knowledge transfer, source task selection, long short-term memory, mutation, crossover.

## I. INTRODUCTION

Multi-task optimization (MTO) [1]–[3], a recently emerging research area in the field of optimization, mainly focuses on investigating how to solve multiple optimization problems at the same time so that the processes of solving relevant problems may help each other via knowledge transfer to boost the overall performance of solving all problems. MTO assumes some useful common knowledge exists for solving related tasks so that the helpful information acquired from addressing one task may be used to help solve another task if these two tasks have certain relatedness [4]. Given its superior performance, MTO has been successfully applied to solve the benchmark optimization problems [5]–[7] and real-world applications [8]–[10]. The research challenges arising from MTO include how to find the helpful source tasks for a target task and how the knowledge from selected source tasks can be extracted, transferred, and reused in a target task.

Evolutionary MTO (EMTO) [11], [12] leverages evolutionary algorithms (EAs) [13] as the optimizer, aiming to unleash the potential of the implicit parallelism featured in EAs for solving MTO problems, where multiple optimization problems are addressed within a unified search space and knowledge is typically represented in the form of promising solutions and transferred via certain evolutionary operations such as crossover and mutation. The development of EMTO includes multifactorial evolutionary algorithm (MFEA) [11] that is one of the most representative EMTO built on the genetic algorithm (GA), multitasking coevolutionary particle swarm optimization (MT-CPSO) that employs multiple swarms for solving multiple tasks [14], an adaptive evolutionary multi-task optimization (AEMTO) framework that can adaptively choose the source tasks with probabilities for each target task working with differential evolution (DE) [15], an evolutionary multitasking-based constrained multi-objective optimization (EMCMO) framework developed to solve constrained multi-objective optimization problems by incorporating GA [6], etc., from which different EAs are involved and their advantages are adopted to exchange knowledge among different tasks.

EMTO has been applied to address regression and classification problems [16], [17]. A co-evolutionary multitasking learning (MTL) approach was proposed in [18] to solve a tropical cyclone wind-intensity prediction problem, where a multi-step ahead prediction problem is formulated as multiple one-step ahead prediction tasks with knowledge represented as a certain part of the neural network. A binary version of an existing multitasking multi-swarm optimization was proposed in [8] to find the optimal feature subspace for each base learner in an ensemble classification model. In [19], an evolutionary multitasking (EMT) ensemble learning model was proposed to solve the hyperspectral image classification problem by modeling feature selection (FS) as an MTO problem. An EMT-based FS method named PSO-EMT was proposed in [20] for solving the high-dimensional classification problem. PSO-EMT mainly focuses on converting a high-dimensional FS

problem into several low-dimensional FS tasks and solving these tasks while enabling knowledge transfer between them.

In this paper, we propose a multi-task optimization based co-training (MTO-CT) framework which trains multiple prediction models simultaneously, where an inter-task knowledge transfer module is designed to transfer and reuse knowledge (represented as model parameters) between these training tasks to facilitate solving them. The long short-term memory (LSTM) [21] based model is employed as the predictor and optimized by a gradient descent (GD) based optimization method for all tasks. The predictor for each task has the same structure. In the inter-task knowledge transfer module, to decide which source tasks to be selected and the amount of knowledge within them to be transferred to help solve the target task, probability-based source task selection [15] is applied, where probability matching (PM) [22] is used to calculate the selection probabilities of all source tasks w.r.t. the current target task, and then stochastic universal selection (SUS) [23] is applied to select the most helpful ones from all sources tasks. Evolutionary operations are then applied to reuse the knowledge from the selected source tasks in the target task. Since this paper is to verify the superiority of MTO in addressing multiple tasks simultaneously, the proposed MTO-CT is compared with the single-task prediction (STP) model without knowledge transfer, i.e., solving each task in a standalone way, under the same settings.

We use electricity consumption data from five states in Australia, i.e., VIC, NSW, SA, QLD, and TAS, to create two sets of tasks at different scales: a) one-step ahead prediction over five states (five tasks) and b) 6-step, 12-step, 18-step, and 24-step ahead prediction for each state (20 tasks), where electricity consumption data in different states share some common patterns. Also, in the multi-step ahead prediction problem, the next-step prediction depends on the knowledge of the previously predicted steps, which is an implicit form of common knowledge across different prediction tasks and makes it reasonable to regard prediction at different steps ahead as related tasks. In comparison to STP, the results on these two sets of tasks verify the superiority of MTO-CT.

The rest of this paper is organized as follows. Section II describes the problem formulated and the background knowledge. The proposed method and its implementation are presented in Section III. Section IV reports and discusses experiments. Conclusions and some planned future work are given in Section V.

## II. PROBLEM DEFINITION AND BACKGROUND

This section will firstly introduce the problem defined. Then the background of LSTM is presented.

### A. Problem Definition

Suppose there are $m$ time series $\mathbf{X} = \{\mathbf{x}_1, ..., \mathbf{x}_m\}$, $\mathbf{x}_i = \{x_{i,1}, ..., x_{i,l_i}\}$, $i \in \{1, ..., m\}$, where $l_i$ is the length of the $i^{th}$ time series. For any time series $i \in \{1, ..., m\}$, there are $p$ different prediction purposes (e.g., different steps ahead prediction). An MTO-CT problem is defined as solving

$n = mp$ prediction tasks at the same time. Given a predictor $h(\cdot)$, any prediction task $j \in \{1, ..., n\}$ can be defined by $h_j(\widetilde{\mathbf{x}}_j; \mathcal{P}_j) \rightarrow \hat{\mathbf{y}}_j$, where $\mathcal{P}_j$ denotes the parameter set of $h_j(\cdot)$ and $(\widetilde{\mathbf{X}}, \mathbf{Y}) = \{(\widetilde{\mathbf{x}}_1, \mathbf{y}_1), (\widetilde{\mathbf{x}}_2, \mathbf{y}_2), ..., (\widetilde{\mathbf{x}}_n, \mathbf{y}_n)\}$ represents the training set for all $n$ task.

Since the target task $j$ may benefit from addressing a source task $k \in \{1, ..., n\}, k \neq j$ via knowledge transfer, knowledge from the source task (i.e., $\mathcal{P}_k \in \{\mathcal{P}_1, \mathcal{P}_2..., \mathcal{P}_n\}, \mathcal{P}_k \neq \mathcal{P}_j$) can be used to help boosting the prediction performance of the $j^{th}$ task. During the update of $\mathcal{P}_j$, knowledge from some selected source tasks based on certain probabilities according to their historical helpfulness is extracted, transferred, and reused to generate $\mathcal{P}_j^{new}$ for the $j^{th}$ target task to help improving its prediction performance.

### B. Long Short-Term Memory

Long short-term memory (LSTM), as a special kind of recurrent neural network (RNN), was proposed in 1997 [21] to overcome the shortcomings of recurrent backpropagation for learning to store information over extended time intervals. LSTM is explicitly designed to avoid the long-term dependency problem and remember information for long periods of time. Similar to the general RNNs, LSTM has a chain of repeating cells of an NN. The structure of an LSTM with one cell is illustrated as Fig. 1, from which we can see there are a cell state ($C_{t-1}$) and three gates, i.e., forget gate ($f_t$), input gate ($i_t$), and output gate ($O_t$).
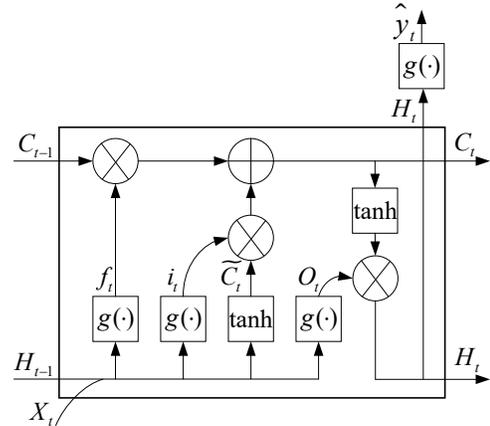


Fig. 1. The structure of long short-term memory.

Given the inputs of current timestamp $X_t$, the hidden state and the cell state of the previous timestamp $H_{t-1}$ and $C_{t-1}$, three gates $f_t$, $i_t$, and $O_t$, and the next cell state $C_t$ can be obtained as follows:

$$
\begin{aligned}
f_t &= g(W_f[H_{t-1}, X_t] + b_f]) \\
i_t &= g(W_i[H_{t-1}, X_t] + b_i]) \\
\tilde{C}_t &= \tanh(W_c[H_{t-1}, X_t] + b_c]) \\
O_t &= g(W_o[H_{t-1}, X_t] + b_o]) \\
C_t &= C_{t-1}f_t + i_t\tilde{C}_t \\
H_t &= O_t \tanh(C_t)
\end{aligned}
\tag{1}
$$

In (1), $g(\cdot)$ in three gates is sigmoid function. With the current hidden state $H_t$, the prediction value $\hat{y}_t$ can be calculated according to:

$$\hat{y}_t = g(W_y H_t + b_y) \tag{2}$$

The activation function $g(\cdot)$ in (2) is sigmoid function in regression problems. The weights $W_f, W_i, W_c, W_o, W_y$ and biases $b_f, b_i, b_c, b_o, b_y$ over different cells are same. To obtain the optimal prediction result is to obtain the optimal weights and biases.

$$\min \frac{1}{NT} \sum_{t=1}^{T} \sum_{s=1}^{N} L(y_{t,s}, \hat{y}_{t,s}) \tag{3}$$
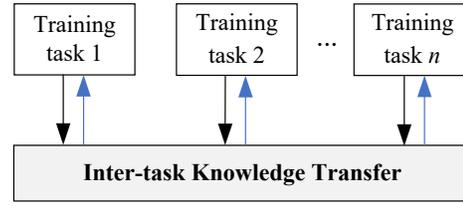
The parameter set $\mathcal{P} = \{W_f, W_i, W_c, W_o, W_y, b_f, b_i, b_c, b_o, b_y\}$ can be learned via any suitable optimization method using the loss function in (3), where the real values are denoted as $\mathbf{y}_t = \{y_{t,1}, y_{t,2}, ..., y_{t,N}\}$, $\hat{\mathbf{y}}_t = \{\hat{y}_{t,1}, \hat{y}_{t,2}, ..., \hat{y}_{t,N}\}$ are the predicted values, $L(\cdot)$ is the evaluation method, $N$ represents the number of samples, and $T$ denotes the number of steps to be predicted.
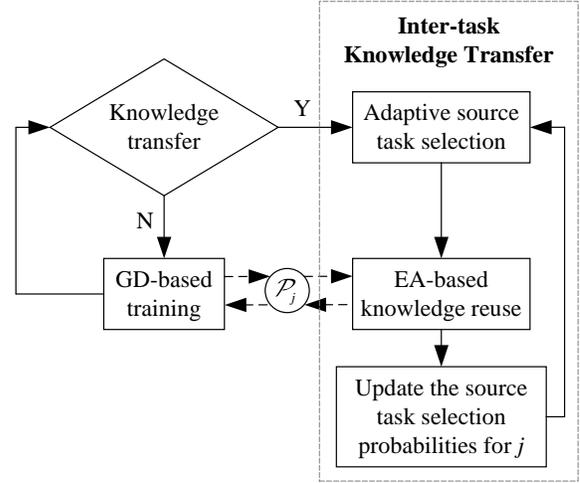
## III. THE PROPOSED METHOD

We will first describe the proposed MTO-CT framework and then elaborate its inter-task knowledge transfer module responsible for selecting the most helpful source tasks in a probabilistic manner, adapting task selection probabilities, and reusing the knowledge from the selected source tasks to assist in the target task. We will also introduce an implementation of the MTO-CT framework.

### A. Framework

The proposed MTO-CT framework is illustrated in Fig. 2, where Fig. 2(a) shows the diagram of the co-training process with $n$ different tasks and Fig. 2(b) describes the individual training process for each task $j$, $j \in \{1, ..., n\}$. As shown in Fig. 2(a), all tasks are solved iteratively. In each iteration, each task is addressed independently with GD-based training before inter-task knowledge transfer is applied. After that, as illustrated in Fig. 2(b), if the knowledge transfer condition satisfies, e.g., for the $j^{th}$ task, the inter-task knowledge transfer will be applied. It first makes the adaptive source task selection, which consists of calculating the selection probabilities of $n-1$ source tasks according to their historical helpfulness in improving the performance of the $j^{th}$ task and selecting the candidates from these source tasks to extract knowledge. EA-based knowledge reusing mainly uses the operations in the EA to create the knowledge to be transferred and reused. The newly generated knowledge is utilized via updating $\mathcal{P}_j$. Finally, the effectiveness of the selected source tasks is quantified and used to update their selection probabilities for the next iteration. Notably, the MTO-CT framework can be treated as a special instance of the AEMTO framework [15], where training instead of general optimization is incorporated.



Fig. 2. The illustration of MTO-CT framework: (a) the overall diagram and (b) the individual training process for each task $j$, $j \in \{1, ..., n\}$.

### B. Inter-task Knowledge Transfer

The inter-task knowledge transfer module consists of choosing the most helpful source tasks to help a target task based on their selection probabilities, transferring and reusing the extracted knowledge from the selected source tasks in the target task, and updating the selection probability of each source based on their helping performance.

*1) Source Task Selection:* For any task $j$, $j \in \{1, ..., n\}$, it has $n-1$ source tasks. In the inter-task knowledge transfer module in Fig. 2(b), the first step is to decide which source task(s) to be selected from the $n-1$ candidates. The source tasks that are more historically helpful may provide more useful knowledge. We calculate the probability of any source task $i$ according to its historical success rate in helping the target task $j$ iteration by iteration. We use $\mathbf{q}_j = \{q_j^i | i = 1, ..., n, i \neq j\}$ obtained from (6) to denote the estimated helpfulness of each source task to the $j^{th}$ target task. With the obtained probabilities, the next step is to select $n_s$ source tasks from all $n-1$ candidates. We use SUS [15], [23] with the source task probabilities $\mathbf{q}_j$ to select $n_s$ source tasks for the target task $j$, $j \in \{1, ..., n\}$.

*2) Knowledge Transfer and Reuse:* After selecting $n_s$ source tasks via SUS, it is important to determine the amount of knowledge to be extracted from each of them and transferred to the target task $j$, $j \in \{1, ..., n\}$, given that the source task with larger probability may provide more helpful

knowledge to help solving the target task. We use the mutation operation used in differential evolution (DE) [24] to generate a mutant $\mathcal{P}_j^{new}$ based on $n_s$ selected source tasks. In this work, we set $n_s = 3$ and adopt a popular DE mutation strategy "DE/rand/1" to produce a mutant as follows:

$$\mathcal{P}_j^{new} = \mathcal{P}_{j_1} + F \cdot (\mathcal{P}_{j_2} - \mathcal{P}_{j_3}) \tag{4}$$

where $j_1, j_2, j_3 \in [1, n], j_1 \neq j_2 \neq j_3 \neq j$ denote three integers yielded via SUS. $F \in [0, 1]$ is a positive real-valued control parameter for scaling the difference vector.

To reuse the knowledge from the selected source tasks in the target task $j$, we apply the binomial crossover operation used in DE to the generated mutant and the target $\mathcal{P}_j$ to generate a new candidate solution as follows:

$$\mathcal{P}_j^{new,d} = \begin{cases} \mathcal{P}_j^{new,d} & \text{if } rand_d[0,1] \leq CR \\ \mathcal{P}_j^d & \text{otherwise} \end{cases} \tag{5}$$

where $d \in \{1, \ldots, D\}$ and $D$ denotes the number of elements in $\mathcal{P}_j$, $j \in \{1, \ldots, n\}$ and $CR \in [0, 1]$ denotes the real-valued crossover rate. $\mathcal{P}_j^{new}$ and $\mathcal{P}_j$ will then compete for survival.

*3) Source Task Selection Probability Update:* The selection probability of each source task is initialized to a very small positive value. In each iteration, after reusing the knowledge from the $i^{th}$ source task in the $j^{th}$ target task, the corresponding helpfulness in the current iteration is measured via the reward $r_j^i$, which is then applied to update $q_j^i$ according to:

$$q_j^i = \alpha q_j^i + (1 - \alpha) r_j^i \tag{6}$$

$r_j^i$ in this work is calculated by the successful rate of the $i^{th}$ task helping the $j^{th}$ task, i.e., $r_j^i = ns_j^i / (na_j^i + \epsilon)$, where $na_j^i$ and $ns_j^i$ denote the total number of times for the $i^{th}$ task selected to help the $j^{th}$ task over a certain period of time and the times that this help leads to the newly generated candidate solution to replace the target one. $\epsilon$ is a quite small positive value to avoid the issue of division by zero.

*C. Implementation*

We implement the MTO-CT framework by using an LSTM-based prediction model for solving each of $n$ tasks. Given only a single time series is considered in this work, we adopt a less typical way to formulate the LSTM-based prediction task. Specifically, the input is defined as the time series values in a time window of $n_f$ consecutive timestamps and the output is defined as the time series values for $1, \ldots, T_j$ steps ahead immediately following this window. Each LSTM cell has a single hidden layer and takes as inputs all $n_f$ time series values in the window as well as the hidden and cell states, where the first cell outputs the one-step ahead prediction, the second cell outputs the two-step ahead prediction and so on till the required $T_j$-step ahead prediction for the $j^{th}$ task is generated. As such, the total number of LSTM cells used is equivalent to $T_j$. The number of hidden neurons in a cell is set to $n_h$. This is different from a more typical way to formulating the LSTM-based prediction task in a "many-to-many" manner, where each LSTM cell is fed in with only

---

**Algorithm 1:** Implementation of MTO-CT

**Input:** $(\widetilde{\mathbf{X}}, \mathbf{Y}) = \{(\widetilde{\mathbf{x}}_1, \mathbf{y}_1), (\widetilde{\mathbf{x}}_2, \mathbf{y}_2), \ldots, (\widetilde{\mathbf{x}}_n, \mathbf{y}_n)\}$, *MaxIter*,
$\quad CR = 0.5, F = 0.2, n_s = 3, r_j^i = 0, ns_j^i = 0,$
$\quad na_j^i = 0, q_j^i = 0.005, j \in \{1, \ldots, n\},$
$\quad j \in \{1, \ldots, n\}, i \neq j, \mathcal{T} = \{T_1, T_2, \ldots, T_n\},$
$\quad \alpha = 0.3, N, D, \#Iter = 0$

1 **for** $j \to 1 : n$ **do**
2 $\quad$ Initialize the parameter set $\mathcal{P}_j$ from the standard normal distribution
3 **end**
4 **while** *#Iter < MaxIter* **do**
5 $\quad$ **for** $j \to 1 : n$ **do**
6 $\quad\quad$ Evaluate the parameter set $\mathcal{P}_j$ on the $j^{th}$ task using (7), denoted as $L_j$
$\quad\quad$ // Inter-task knowledge transfer starts
7 $\quad\quad$ **for** $i \to 1 : n$ & $i \neq j$ **do**
8 $\quad\quad\quad$ Calculate each source task selection probability according to (6) to obtain the updated $q_j^i$
9 $\quad\quad$ **end**
10 $\quad\quad$ Perform SUS operation [23] to select $n_s$ source tasks, i.e., $j_1, j_2, \ldots, j_{n_s}, k \in \{1, \ldots, n_s\}, j_k \in \{1, \ldots, n\}, j_k \neq j$
11 $\quad\quad$ Perform mutation operation according to (4) with the selected source tasks $j_1, j_2, \ldots, j_{n_s}$ to obtain the $\mathcal{P}_j^{new}$
12 $\quad\quad$ **for** $d \to 1 : D$ **do**
13 $\quad\quad\quad$ Perform crossover operation with (5) to update $\mathcal{P}_j^{new,d}$
14 $\quad\quad$ **end**
15 $\quad\quad$ Evaluate the newly generated $\mathcal{P}_j^{new}$ with (7) to obtain the performance $L_j^{new}$
16 $\quad\quad$ **if** $L_j^{new} < L_j$ **then**
17 $\quad\quad\quad$ $\mathcal{P}_j^{new} \to \mathcal{P}_j$
18 $\quad\quad\quad$ $L_j^{new} \to L_j$
19 $\quad\quad\quad$ $ns_j^{j_k} = ns_j^{j_k} + 1$
20 $\quad\quad$ **end**
21 $\quad\quad$ $na_j^{j_k} = na_j^{j_k} + 1, k \in \{1, \ldots, n_s\}, j_k \in \{1, \ldots, n\}$
$\quad\quad$ // Inter-task knowledge transfer ends
22 $\quad\quad$ Update the parameter set $\mathcal{P}_j$ with Adam algorithm [25]
23 $\quad$ **end**
24 $\quad$ $\#Iter = \#Iter + 1$
25 **end**
**Output:** $\mathcal{P}_1^*, \mathcal{P}_2^*, \ldots, \mathcal{P}_n^*, \hat{\mathbf{y}}_1, \hat{\mathbf{y}}_2, \ldots, \hat{\mathbf{y}}_n$

---

one time series value at a certain timestamp. For each LSTM-based model, the parameters to be optimized (trained) include $\{W_f^j, W_i^j, W_c^j, W_o^j, W_y^j, b_f^j, b_i^j, b_c^j, b_o^j, b_y^j\}$ encoded via $\mathcal{P}_j$ of dimension size $D$.

Each prediction task $j$, $j \in \{1, \ldots, n\}$ is addressed by an adaptive moment estimation (Adam) [26], which is a first-order GD-based optimization method with the adaptive estimates of lower-order moments [25]. The inter-task knowledge transfer module is performed to explore if the information extracted from $n_s$ source tasks can further boost the prediction accuracy of the target task $j$.

We employ the root mean square error (RMSE) to define the

loss function in (3) for any task $j, j \in \{1, \ldots, n\}$ as follows:

$$\min \sqrt{\frac{1}{NT_j} \sum_{t=1}^{T_j} \sum_{s=1}^{N} (y_j^{t,s} - \hat{y}_j^{t,s})^2} \qquad (7)$$

where $T_j, j \in \{1, \ldots, n\}$ denotes the time steps ahead to be predicted in the $j^{th}$ task, $\hat{\mathbf{y}}_j$ denotes the predicted result w.r.t. the ground truth $\mathbf{y}_j$, where $\hat{\mathbf{y}}_j = \{\hat{y}_j^1, \hat{y}_j^2, ..., \hat{y}_j^{T_j}\}, \hat{y}_j^t \in \Re^{1 \times N}$ and $\hat{y}_j^t = \{y_j^{t,1}, y_j^{t,2}, ..., y_j^{t,N}\}, t \in \{1, \ldots, T_j\}$.

In each iteration, $\mathcal{P}_j$ is replaced by $\mathcal{P}_j^{new}$ via the inter-task knowledge transfer when $\mathcal{P}_j$ has worse performance than $\mathcal{P}_j^{new}$ in terms of (7). This repeats until the maximum number of iterations (*MaxIter*) is reached. The implementation of MTO-CT is detailed in Algorithm 1.

## IV. RESULTS

We will first present the data information and experimental settings. Then the prediction performance on these two different sets of tasks are presented and compared with STP to demonstrate the superiority of MTO-CT.

### A. Data Description and Experimental Settings

The data is downloaded from Australian Energy Market Operator (AEMO)[1]. It includes electricity consumption data collected at 30-minute intervals from 01 November 2020 to 30 November 2021 for five states (VIC, NSW, SA, QLD, TAS) in Australia. We create two sets of tasks at different scales: (1) Set A: one-step ahead prediction across five states (five prediction tasks); (2) Set B: multi-step ahead (e.g., 6, 12, 18, 24) prediction for each of these five states (20 prediction tasks). For both sets, the time windows used as inputs are set as 24 (i.e., using the first 12 hours to predict the next step, next several steps, or the rest of the same day). For each of the tasks, there are 395 samples in total, where training and testing samples occupy 80% (316) and 20% (79), respectively. For each state, the data is normalized to $[0, 1]$ using the min-max normalization.

The aim of this paper is to investigate if MTO can help improve the prediction accuracy when having multiple prediction tasks to be addressed simultaneously. We compare the results of MTO-CT with that of STP, which addresses every single task independently without inter-task knowledge transfer. The number of hidden neurons in LSTM is $n_h = 10$. In GD-based optimization method, i.e., Adam in this paper, learning rate is $l_r = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 1e - 8$. To guarantee the comparison under the same number of evaluations, *MaxIter* = 20000 in STP and *MaxIter* = 10000 in MTO-CT considering the inter-task knowledge transfer operation in each iteration. All tasks are independently run ten times.

### B. Results

The performance of MTO-CT is comprehensively studied by comparing it with STP over two sets of tasks, i.e., five and 20 tasks. The training and testing performance (RMSE) of

[1]http://www.nemweb.com.au/REPORTS/Archive/HistDemand/

TABLE I
THE TASK REPRESENTATIONS BY THE NUMBERS.

| States | | VIC | NSW | SA | QLD | TAS |
|---|---|---|---|---|---|---|
| Set A | one-step | 1 | 2 | 3 | 4 | 5 |
| Set B | 6-step | 1 | 2 | 3 | 4 | 5 |
| | 12-step | 6 | 7 | 8 | 9 | 10 |
| | 18-step | 11 | 12 | 13 | 14 | 15 |
| | 24-step | 16 | 17 | 18 | 19 | 20 |

these two sets of tasks over MTO-CT and STP is summarized and discussed, where the results is based on the normalized data. Wilcoxon signed-rank test at the 0.05 level is performed to estimate the significance of the difference between MTO-CT and STP. The better performance over the statistical test is labeled in bold. We use '+', '=', and '-' to indicate that the respective model has better, same, and worse performance than the other(s). To better understand what the task number represents in the following results, Table. I gives the task representations with numbers.

TABLE II
AVERAGE RMSE OVER TRAINING AND TESTING SETS FOR ONE-STEP AHEAD PREDICTION OVER FIVE STATES ('+', '=', '-' : BETTER, SAME, WORSE).

| Tasks | Training RMSE | | Testing RMSE | |
|---|---|---|---|---|
| | STP | MTO-CT | STP | MTO-CT |
| 1 | 0.05885 | **0.05747** | 0.07012 | **0.06824** |
| 2 | 0.0528 | **0.04966** | 0.0592 | **0.05427** |
| 3 | 0.05237 | **0.05141** | 0.05283 | **0.0502** |
| 4 | 0.06332 | **0.05859** | 0.0588 | **0.05563** |
| 5 | 0.05721 | **0.05457** | 0.05918 | **0.05632** |
| +/=/- | 0/0/5 | 5/0/0 | 0/0/5 | 5/0/0 |

### C. Results of Set A

Table. II reports the training and testing performance evaluated by mean RMSE over ten independent runs for MTO-CT and STP (without inter-task knowledge transfer), where the best mean RMSE for each task is labeled bold if it is significantly better with the statistical test. The result is based on one-step ahead prediction over VIC, NSW, SA, QLD, and TAS, respectively. By comparing the mean RMSE of MTO-CT and STP, it is obvious that MTO-CT outperforms both training and testing sets over all tasks from the labeled bold values and the total number of tasks it wins. This shows that the helpful knowledge reuse of the selected source tasks leads to significant improvement in the performance of the target task so that the accuracy of all tasks can be enhanced, which also demonstrates the effectiveness of inter-task knowledge transfer in MTO-CT.

Fig. 3 illustrates the training and testing RMSE distribution of ten independent runs over each task for STP and MTO-
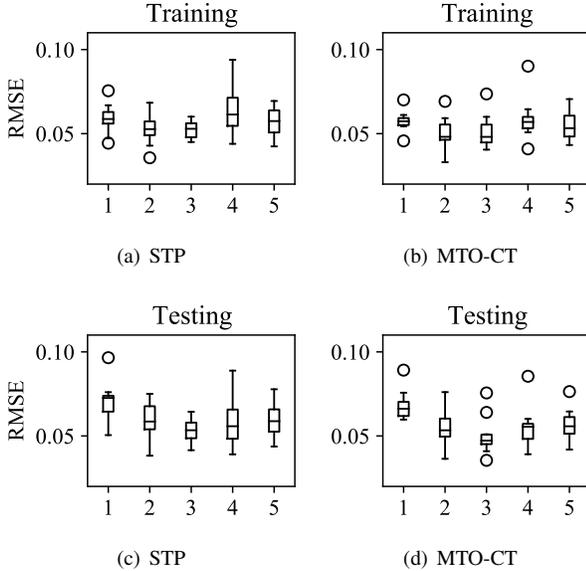
Fig. 3. Box plots of training and testing RMSE on five tasks under STP and MTO-CT over ten independent runs.

CT. Fig. 3(a) and Fig. 3(b) show that STP and MTO-CT have slight difference over the training performance from the view of distribution, but we still can see that the training RMSE on MTO-CT is better than STP. Even though MTO-CT has more outliers, fewer results deviate from the average value for each of the tasks. Similarly, comparing Fig. 3(c) and Fig. 3(d), the results on STP have more differences, resulting in higher average RMSE on each task, verified by the results in Table. II.

### D. Results of Set B

Table. III shows the average RMSE across ten independent runs over training and testing sets for STP and MTO-CT. The results are from 20 tasks that include 6-step, 12-step, 18-step, and 24-step ahead prediction over five states (the representations of the task numbers can be found in Table. I). From Table. III, we can see some tasks over MTO-CT have worse training performance than STP such as tasks 8 (12-step ahead prediction in SA), 9 (12-step ahead prediction in QLD), 14 (18-step ahead prediction in QLD), 17 (24-step ahead prediction in NSW), and 18 (24-step ahead prediction in SA). MTO-CT and STP have the same significant level on task 3 (6-step ahead prediction in SA) on the training set from the result of the statistical test, even though the average RMSE is slightly different. Among all 20 tasks, MTO-CT outperforms STP on 14 tasks on the training set. For the testing set, MTO-CT leads to better performance on 17 tasks than STP, where tasks 6 (12-step ahead prediction in VIC), 9, and 18 have worse performance (higher mean RMSE) with MTO-CT. The result of addressing 20 tasks simultaneously with inter-task knowledge transfer further demonstrates the superiority of MTO-CT, given that it outperforms STP on 14 and 17 tasks for training and testing sets, respectively.

TABLE III
AVERAGE RMSE OVER TRAINING AND TESTING SETS FOR FIVE STATES
ACROSS 6-STEP, 12-STEP, 18-STEP, AND 24-STEP AHEAD PREDICTION
('+', '=', '-' : BETTER, SAME, WORSE).

|  | Training RMSE | | Testing RMSE | |
|---|---|---|---|---|
| Tasks | STP | MTO-CT | STP | MTO-CT |
| 1 | 0.09486 | **0.09179** | 0.10158 | **0.09783** |
| 2 | 0.0903 | **0.0879** | 0.09211 | **0.09006** |
| 3 | **0.08879** | 0.08655 | 0.08504 | **0.08016** |
| 4 | 0.09737 | **0.09588** | 0.09111 | **0.08963** |
| 5 | 0.07772 | **0.07625** | 0.08044 | **0.07921** |
| 6 | 0.12331 | **0.12401** | **0.12088** | 0.12204 |
| 7 | 0.11972 | **0.11756** | 0.12252 | **0.11975** |
| 8 | **0.11451** | 0.11472 | 0.11011 | **0.10907** |
| 9 | **0.13277** | 0.13398 | **0.12709** | 0.1276 |
| 10 | 0.09543 | **0.09385** | 0.09304 | **0.09074** |
| 11 | 0.12747 | **0.12587** | 0.1283 | **0.12392** |
| 12 | 0.12901 | **0.12548** | 0.13309 | **0.12913** |
| 13 | 0.12717 | **0.12657** | 0.12319 | **0.12159** |
| 14 | **0.13083** | 0.13175 | 0.12711 | **0.12697** |
| 15 | 0.10001 | **0.09781** | 0.09744 | **0.09531** |
| 16 | 0.1185 | **0.1177** | 0.11515 | **0.11377** |
| 17 | **0.11527** | 0.1154 | 0.11907 | **0.119** |
| 18 | **0.11751** | 0.12034 | **0.11562** | 0.11859 |
| 19 | 0.13719 | **0.13468** | 0.13418 | **0.13309** |
| 20 | 0.10684 | **0.10257** | 0.10331 | **0.10123** |
| +/=/- | 5/1/14 | 14/1/5 | 3/0/17 | 17/0/3 |

Fig. 4 illustrates the distribution of the training and testing RMSE on 20 tasks under STP and MTO-CT. For most of the tasks, Fig. 4(a) and Fig. 4(b) show similar distribution over ten runs, except for task 18, which is significantly different and also has worse performance on MTO-CT. For the testing RMSE as presented in Fig. 4(c) and Fig. 4(d), task 18 on MTO-CT cannot compete with STP as well. However, for most of the rest, MTO-CT outperforms STP, further verified by Table. III. Therefore, knowledge transfer among tasks in MTO-CT leads to better prediction performance for most tasks.

### V. CONCLUSIONS AND FUTURE WORK

We proposed an MTO-CT framework to solve multiple prediction tasks simultaneously, where an inter-task knowledge transfer module is designed to transfer and share knowledge among different tasks so that the overall performance of solving each task can be improved. MTO-CT employs an LSTM based model as the predictor and represents the knowledge as the connection weights and biases in LSTM. The inter-task knowledge transfer module is responsible for selecting the source tasks (w.r.t. a target task) from which the knowledge is extracted, extracting the knowledge, and reusing the extracted

Training

(a) STP

Training

(b) MTO-CT
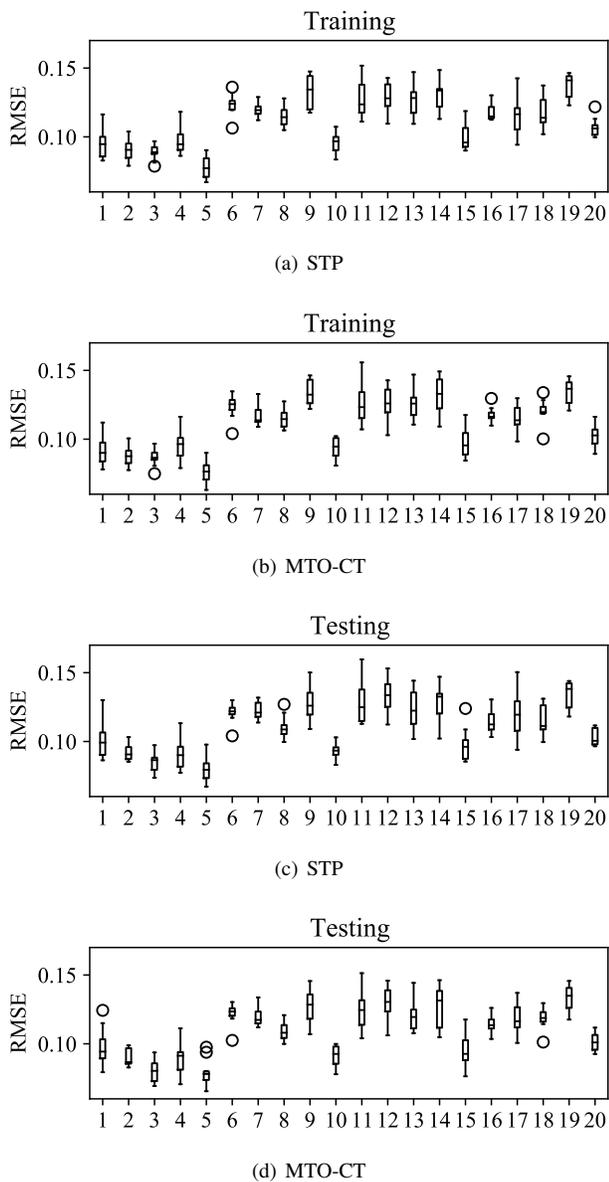
Testing

(c) STP

Testing

(d) MTO-CT

Fig. 4. Box plots of training and testing RMSE on 20 tasks under STP and MTO-CT over ten independent runs.

knowledge in the target task. The performance of MTO-CT is tested on two sets of tasks at different scales, i.e., five tasks and 20 tasks. The superiority of MTO-CT in terms of prediction accuracy is demonstrated in comparison to STP which solves each task in a standalone way without inter-task knowledge transfer. Our future work includes enriching the input by incorporating additional time series data like temperature, evaluating the performance of MTO-CT for co-training more LSTM variants or other types of prediction models [27], and applying MTO-CT to other applications that we worked on previously like graph matching [28], feature extraction [29] service composition [30].

REFERENCES

[1] H. Song, A. K. Qin, P. W. Tsai, and J. J. Liang, "Multitasking multi-swarm optimization," in *Proc. of 2019 IEEE Congress on Evolutionary Computation (CEC)*. IEEE, 2019, pp. 1937–1944.
[2] K. Swersky, J. Snoek, and R. P. Adams, "Multi-task bayesian optimization," *Advances in neural information processing systems*, vol. 26, 2013.
[3] Q. Chen, X. Ma, Y. Yu, Y. Sun, and Z. Zhu, "Multi-objective evolutionary multi-tasking algorithm using cross-dimensional and prediction-based knowledge transfer," *Information Sciences*, vol. 586, pp. 540–562, 2022.
[4] Z. Liang, J. Zhang, L. Feng, and Z. Zhu, "A hybrid of genetic transform and hyper-rectangle search strategies for evolutionary multi-tasking," *Expert Systems with Applications*, vol. 138, p. 112798, 2019.
[5] A. Gupta, Y. S. Ong, and L. Feng, "Insights on transfer optimization: Because experience is the best teacher," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 2, no. 1, pp. 51–64, 2017.
[6] K. Qiao, K. Yu, B. Qu, J. Liang, H. Song, and C. Yue, "An evolutionary multitasking optimization framework for constrained multi-objective optimization problems," *IEEE Transactions on Evolutionary Computation*, 2022.
[7] K. Qiao, K. Yu, B. Qu, J. Liang, H. Song, C. Yue, H. Lin, and K. C. Tan, "Dynamic auxiliary task-based evolutionary multitasking for constrained multi-objective optimization," *IEEE Transactions on Evolutionary Computation*, 2022.
[8] B. Zhang, A. K. Qin, and T. Sellis, "Evolutionary feature subspaces generation for ensemble classification," in *Proc. of The Genetic and Evolutionary Computation Conference*, 2018, pp. 577–584.
[9] R. Ye and Q. Dai, "Multitl-kelm: A multi-task learning algorithm for multi-step-ahead time series prediction," *Applied Soft Computing*, vol. 79, pp. 227–253, 2019.
[10] Z. Wu, Q. Li, and X. Xia, "Multi-timescale forecast of solar irradiance based on multi-task learning and echo state network approaches," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 1, pp. 300–310, 2020.
[11] A. Gupta, Y. S. Ong, and L. Feng, "Multifactorial evolution: toward evolutionary multitasking," *IEEE Transactions on Evolutionary Computation*, vol. 20, no. 3, pp. 343–357, 2015.
[12] D. Liu, S. Huang, and J. Zhong, "Surrogate-assisted multi-tasking memetic algorithm," in *Proc. of 2018 IEEE Congress on Evolutionary Computation (CEC)*. IEEE, 2018, pp. 1–8.
[13] K. A. De Jong, *Evolutionary Computation: A Unified Approach*. Cambridge, MA, USA: MIT Press, 2016.
[14] M. Y. Cheng, A. Gupta, Y. S. Ong, and Z. W. Ni, "Coevolutionary multitasking for concurrent global optimization: With case studies in complex engineering design," *Engineering Applications of Artificial Intelligence*, vol. 64, pp. 13–24, 2017.
[15] H. Xu, A. Qin, and S. Xia, "Evolutionary multi-task optimization with adaptive knowledge transfer," *IEEE Transactions on Evolutionary Computation*, 2021.
[16] R. Chandra, A. Gupta, Y. S. Ong, and C. K. Goh, "Evolutionary multi-task learning for modular knowledge representation in neural networks," *Neural Processing Letters*, vol. 47, no. 3, pp. 993–1009, 2018.
[17] R. Chandra, Y. S. Ong, and C. K. Goh, "Co-evolutionary multi-task learning for dynamic time series prediction," *Applied Soft Computing*, vol. 70, pp. 576–589, 2018.
[18] R. Chandra, "Dynamic cyclone wind-intensity prediction using co-evolutionary multi-task learning," in *Proc. of International Conference on Neural Information Processing*. Springer, 2017, pp. 618–627.
[19] J. Shi, T. Shao, X. Liu, X. Zhang, Z. Zhang, and Y. Lei, "Evolutionary multitask ensemble learning model for hyperspectral image classification," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 936–950, 2020.
[20] K. Chen, B. Xue, M. Zhang, and F. Zhou, "Evolutionary multitasking for feature selection in high-dimensional classification via particle swarm optimisation," *IEEE Transactions on Evolutionary Computation*, 2021.
[21] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
[22] D. Thierens, "An adaptive pursuit strategy for allocating operator probabilities," in *Proc. of The 7th Aannual Conference on Genetic and Evolutionary Computation*, 2005, pp. 1539–1546.
[23] J. E. Baker *et al.*, "Reducing bias and inefficiency in the selection algorithm," in *Proc. of The Second International Conference on Genetic Algorithms*, vol. 206, 1987, pp. 14–21.

[24] A. K. Qin, V. L. Huang, and P. N. Suganthan, "Differential evolution algorithm with strategy adaptation for global numerical optimization," *IEEE Transactions on Evolutionary Computation*, vol. 13, no. 2, pp. 398–417, 2008.

[25] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[26] S. Sun, Z. Cao, H. Zhu, and J. Zhao, "A survey of optimization methods from a machine learning perspective," *IEEE Transactions on Cybernetics*, vol. 50, no. 8, pp. 3668–3681, 2019.

[27] A. K. Qin and P. N. Suganthan, "Initialization insensitive LVQ algorithm based on cost-function adaptation," *Pattern Recognition*, vol. 38, no. 5, pp. 773–776, 2005.

[28] M. Gong, Y. Wu, Q. Cai, W. Ma, A. K. Qin, Z. Wang, and L. Jiao, "Discrete particle swarm optimization for high-order graph matching," *Information Sciences*, vol. 328, pp. 158–171, 2016.

[29] A. K. Qin, P. N. Suganthan, and M. Loog, "Uncorrelated heteroscedastic lda based on the weighted pairwise chernoff criterion," *Pattern Recognition*, vol. 38, no. 4, pp. 613–616, 2005.

[30] S. Mistry, A. Bouguettaya, H. Dong, and A. K. Qin, "Metaheuristic optimization for long-term iaas service composition," *IEEE Transactions on Services Computing*, vol. 11, no. 1, pp. 131–143, 2018.