# ©CUFF-KT: TACKLING LEARNERS' REAL-TIME LEARNING PATTERN ADJUSTMENT VIA TUNING-FREE KNOWLEDGE STATE-GUIDED MODEL UPDATING

Anonymous authors

Paper under double-blind review

# Abstract

Knowledge Tracing (KT) is a core component of Intelligent Tutoring Systems, modeling learners' knowledge state to predict future performance and provide personalized learning support. Current KT models simply assume that training data and test data follow the same distribution. However, this is challenged by the continuous changes in learners' patterns. In reality, learners' patterns change irregularly at different stages (e.g., different semesters) due to factors like cognitive fatigue and external stress. Additionally, there are significant differences in the patterns of learners from various groups (e.g., different classes), influenced by social cognition, resource optimization, etc. We refer to these distribution changes at different stages and from different groups as intra-learner shift and inter-learner shift, respectively—a task introduced, which we refer to as Real-time Learning Pattern Adjustment (RLPA). Existing KT models, when faced with RLPA, lack sufficient adaptability, because they fail to timely account for the dynamic nature of different learners' evolving learning patterns. Current strategies for enhancing adaptability rely on retraining, which leads to significant overfitting and high time cost problem. To address this, we propose Cuff-KT, comprising a controller and a generator. The controller assigns value scores to learners, while the generator generates personalized parameters for selected learners. Cuff-KT adapts to distribution changes fast and flexibly without fine-tuning. Experiments on one classic and two latest datasets demonstrate that Cuff-KT significantly improves current KT models' performance under intra- and inter-learner shifts, with an average relative increase of 7% on AUC, effectively tackling RLPA.<sup>1</sup>

032 033 034

031

000

001

003

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

028

## 1 INTRODUCTION

For nearly a century, researchers have been dedicated to developing Intelligent Tutoring Systems (ITS) (Pressey, 1926; Kamalov et al., 2023; Zhou et al., 2024). *Knowledge Tracing (KT), as a core component of ITS, aims to model learners' knowledge state during their interactions with ITS to predict their performance on future questions* (Corbett & Anderson, 1994), as shown in Figure 1.
Solving the KT problem can help teachers or systems better identify learners who need further attention and recommend personalized learning materials to them (Liu et al., 2021; Abdelrahman et al., 2023; Liu et al., 2023b).

Reviewing the current research on KT (Piech et al., 2015; Shen et al., 2022; Liu et al., 2023a), we can systematize a dominant paradigm: using learners' historical interaction sequences as training data, encoding them into representations with KT models, and then using these representations to predict future interactions in the test data. This



Figure 1: Illustration of the knowledge tracing (KT).

paradigm simply assumes that the training data and test data come from the same distribution. How ever, this assumption is difficult to hold in real-world scenarios, as it ignores the dynamic nature

<sup>053</sup> 

<sup>&</sup>lt;sup>1</sup>Our code and datasets are available at https://anonymous.4open.science/r/Cuff-KT.

of KT. Specifically, the properties of streaming data (*e.g.*, the correct rate distribution) often change
 across different stages or groups (Zhang et al., 2017; Yang et al., 2023), indicating that the sequential
 patterns of learners at different stages or from different groups dynamically vary between historical
 and future interactions. We refer to these distribution changes across different stages and groups as
 intra-learner shift and inter-learner shift, respectively.

Distribution shift caused by varying sequential patterns undermines current KT models, resulting 060 in deteriorated generalization when serving future data. Figure 2 provides empirical evidence of 061 this issue. We first divide the assist15 data into 4 non-overlapping parts by stage and group (see 062 Section 4.3 for the division method), and calculate the KL-divergence w.r.t. correct rate distribution 063 between the first part and the other parts. We then train the DKT (Piech et al., 2015) model on the 064 first part and test it on the remaining parts. Clearly, as the KL-divergence increases across different stages or groups, the model's predictive performance significantly declines. Therefore, it is crucial 065 to enhance the dynamic adaptability of KT models. To this end, we introduce a new task, Real-time 066 Learning Pattern Adjustment (RLPA), to address the inability of existing KT models to effectively 067 handle distribution changes arising from differing learning patterns across various stages or groups. 068

069 To tackle RLPA, a well-070 known generalization technique is to retrain (e.g., 071 fine-tune) the pre-trained 072 KT model based on data 073 from the current stage or 074 group to achieve person-075 alized learning (Houlsby 076 et al., 2019; Zaken et al., 077 2021; Han et al., 2024). Although fine-tuning based 079 approaches are promising, 080 they may not be the opti-081 mal solution due to two key





 challenges: (i) Overfitting. To achieve personalized learning, fine-tuning based approaches often require retraining the model based on very limited samples with rapidly changing distributions, which may lead to overfitting, potentially reducing its ability to generalize to future distributions. (ii) High time cost. Fine-tuning is very time-consuming as it requires extensive gradient computations to update model parameters, which is cumbersome in real-world scenarios where real-time requirements are common. Therefore, fine-tuning based methods must carefully balance the need to adapt to recent data and maintain robustness to achieve generalization. These challenges prompt us to reconsider the design of better solutions to the RLPA in KT.

Towards this end, we propose a novel method to trackle RLPA in KT, called Controllable, tUning-090 free, Fast, and Flexible Knowledge Tracing (Cuff-KT). Unlike fine-tuning-based approaches that 091 produce updated parameters, the core idea of Cuff-KT is to learn a model parameter generator 092 specific to the current stage or group, generating updating personalized parameters for valuable 093 learners (e.g., those showing significant progress or regression), achieving adaptive generalization. 094 Our Cuff-KT consists of two modules: a controller and a generator. When the data distribution of 095 learners changes due to varying learning patterns, the KT model generalizes worse to the cur-096 rent data and tends to make incorrect evaluations. This implies that the benefit of generating 097 parameters is significant, as generated parameters can appropriately model the current data 098 distribution. The controller, while considering the fine-grained distance between knowledge state distributions across various concepts, is also inspired by the Dynamic Assessment Theory (Vygot-099 sky & Cole, 1978) and integrates coarse-grained changes in correct rates, assigning a value score to 100 each learner. The generator generates parameters for learners selected based on the assigned value 101 scores<sup>2</sup> by the controller and enhances adaptability. Specifically, considering the relative relation-102 ship between question difficulty and learner ability (Rasch, 1993; Shen et al., 2022) and inspired by 103 the dual-tower model in recommendation systems (Huang et al., 2013), the generator models ques-104 tions and responses separately, extracts features through a sequential feature extractor, simulates the 105 distribution of real-time samples from the current stage or group to achieve adaptive generaliza-106

<sup>&</sup>lt;sup>2</sup>The larger a learner's value score, the more likely they are to be selected.

tion through our designed state-adaptive attention, and finally reduces the parameter size through
 low-rank decomposition. Notably, our generator can be inserted into into any layer or generate
 parameters for any layer.

### 111 112 Our contributions are summarized as follows:

113

114

115

116

117

118

119

121

122

123

124

125 126 127

128

- We introduce a new task, **RLPA**, which enhances the adaptability of existing KT models in the realm of personalized learning, addressing the challenges arising from distribution shifts caused by varying sequential patterns of learners across different stages or groups.
- We propose **Cuff-KT**, a controllable, tuning-free, fast, and flexible general neural method, which can effectively generate parameters aligned with the current stage or group's learner distribution and insert them into any layer of existing KT models. It is noteworthy that Cuff-KT is model-agnostic.
- We instantiate one classic KT model and two latest state-of-the-art models. Experiments on one classic dataset and two latest datasets demonstrate that our proposed Cuff-KT generally improves current KT models under both intra- and inter-learner shift. Specifically, the AUC metric, which is most commonly used in KT, has relatively increased by 7% on average, proving that Cuff-KT can effectively tackle RLPA in KT.

# 2 RELATED WORK

129 Knowledge tracing (KT), the task of dynamically modeling a learner's knowledge state over time, 130 traces its origins back to the early 1990s, with early notable contributions by Corbett and Ander-131 son (Corbett & Anderson, 1994). However, with the rise of deep learning, KT research has gained significant momentum, leading to the development of more sophisticated and refined models capa-132 ble of capturing the intricate dynamics of learner learning (Piech et al., 2015; Yeung & Yeung, 2018; 133 Shen et al., 2022; Liu et al., 2023a). DKT (Piech et al., 2015) first applies LSTM to KT to model 134 the complex learners' cognitive process, bringing a leap in performance compared to previous KT 135 models (e.g., BKT (Corbett & Anderson, 1994)). Subsequently, various neural architectures (e.g., 136 attention and graphs) begin to be introduced into KT (Pandey & Karypis, 2019; Nakagawa et al., 137 2019; Ghosh et al., 2020; Pandey & Srivastava, 2020). Meanwhile, some training techniques (e.g., 138 adversarial training and contrastive learning) also start to be used in KT research (Guo et al., 2021; 139 Lee et al., 2022a). Recently, incorporating learning-related information has been explored to en-140 hance the predictive capability of KT models. For instance, DIMKT (Shen et al., 2022) improves 141 KT performance by establishing relationship between learners' knowledge states and question dif-142 ficulty levels, while AT-DKT (Liu et al., 2023a) addresses the issues of sparse representation and personalization in DKT by introducing two auxiliary learning tasks: question tagging prediction 143 and individualized prior knowledge prediction. 144

145 However, surprisingly, to our knowledge, there is a lack of attention to adaptability in KT research, 146 which severely affects the generalization of KT models across different distributions. Meanwhile, 147 related studies (Wong et al., 2022; Wong & Ramasamy, 2024) are limited in their applicable scenarios (e.g., continuously increasing learners or concepts) and do not provide a general method to 148 enhance adaptability in KT. Thanks to the well-known fine-tuning based methods, the adaptability 149 in KT has been enhanced to some extent. However, the challenges posed by overfitting and high 150 time cost of fine-tuning based methods make it difficult to be effectively applied in real-world sce-151 narios (Lv et al., 2023b). Even the recently proposed parameter-efficient fine-tuning based methods 152 (e.g., Adapter-based tuning (Houlsby et al., 2019) and Bias-term Fine-tuning (Zaken et al., 2021)) 153 still incur non-negligible time cost and cannot avoid the potential risk of overfitting. Our Cuff-KT, 154 in contrast, updates KT models under dynamic distributions through controllable parameter genera-155 tion, eliminating the need for retraining and providing a new perspective on enhancing adaptability 156 in the KT community.

157 158

# 3 Methodology

159 160

# 161 In this section, we first define the problem of KT and formalize the RLPA task in KT, then introduce our proposed Cuff-KT method.



Figure 3: Overview of proposed Cuff-KT method.

189

190

195 196

197

198 199

200 201

202

207 208

209

3.1 PROBLEM FORMULATIONS

# 180 3.1.1 KNOWLEDGE TRACING

Formally, let S, Q, and C denote the sets of learners, questions, and concepts, respectively. For each learner  $s \in S$ , their interactions are represented  $X^s = \{x\}_{i=1}^k$  at time-step k, where the interaction x is defined as a 4-tuple, *i.e.*,  $x = (q, \{c\}, r, t)$ , where  $q \in Q, \{c\} \subset C, r, t$  represent the question attempted by the learner s, the concepts associated with the question q, the binary variable indicating whether the learner responds to the question correctly (1 for correct, 0 for incorrect), and the timestamp of the learner's response respectively. The goal of KT is to predict  $\hat{r}_{k+1}$  given the learner's historical interactions X and the current question  $q_{k+1}$  at time-step k + 1.

## 3.1.2 RLPA TASK

191 RLPA aims to address two common shift issues (intra- and inter- shifts) in KT to enhance the adapt-192 ability of existing models. An interaction sequence of a learner s can be divided into multiple stages, 193 assuming each stage has a length of L. At time-step u, the representation of the learner's interaction 194 in that stage is  $X_u^s = X_{u:u+L-1}^s$ . Intra-learner shift is defined as: for any time-step  $u \neq v$ ,

$$|d(\chi_u^s, \chi_v^s)| > \delta,\tag{1}$$

where  $\delta$  is a small threshold.  $\chi_u^s$  and  $\chi_v^s$  represent the distributions of  $X_u^s$  and  $X_v^s$  respectively. d is a distance function (*e.g.*, KL divergence). In contrast, inter-learner shift is:

$$|d(\chi_u^s, \chi_u^{s^*})| > \delta, \tag{2}$$

where  $\chi_u^s$  and  $\chi_u^{s^*}$  represent the distributions of learners s and  $s^*$  at time-step u, respectively.

When equations 1 or 2 hold, the goal of RLPA is to adjust the parameters of the existing KT model in real-time so that the predicted distribution  $\hat{\chi}_v^s$  or  $\hat{\chi}_u^{s^*}$  is as close as possible to the actual distribution:

$$\min_{\hat{\chi}_{v}^{s}} \sum_{x} \chi_{v}^{s}(x) \log(\frac{\chi_{v}^{s}(x)}{\hat{\chi}_{v}^{s}(x)}) \text{ or } \min_{\hat{\chi}_{u}^{s*}} \sum_{x} \chi_{u}^{s*}(x) \log(\frac{\chi_{u}^{s*}(x)}{\hat{\chi}_{u}^{s*}(x)}),$$
(3)

where x is a variable in the sample space.

3.2 CUFF-KT

Figure 3 illustrates an overview of our Cuff-KT method, which consists of two modules: (a) Controller identifies learners with valuable parameter update potential, aiming to reduce the cost of
parameter generation. (b) Generator adjusts network parameters for existing KT models at different stages or for different groups, aiming to enhance adaptive generalization. In our setup, the KT
model is decoupled into a static backbone and a dynamic layer. The generator can be inserted into
any layer of the KT model or generate parameters for any layer (dynamic layer). Finally, we introduce the training strategy for Cuff-KT.

# 3.2.1 CONTROLLER

217 218

The controller can identify learners with dramatic changes in their knowledge state distribution (*i.e.*, valuable learners, often showing progress or regression), aiming to reduce the cost of parameter generation. The controller comprehensively considers both fine-grained and coarse-grained changes in the knowledge states of different learners, as described below.

Fine-grained Changes. At time-step k, the KT model models knowledge states States<sup>j</sup> (*i.e.*, proficiency scores ranging from 0 to 1 for |C| concepts over k time steps) for learner s<sup>j</sup> with number j ( $1 \le j \le n$ , where n is the total number of learners) at different time steps. The States<sup>j</sup> is utilized by the controller to measure the fine-grained distance (*e.g.*, KL-divergence) between the knowledge state distributions across various concept at the intermediate time-step k/2 and current time-step k:

231 232

244 245

251

256 257

258

$$\begin{aligned} &\text{States}_{k/2}^{*j} = \text{norm}(\text{States}_{k/2}^{j}), \\ &\text{States}_{k}^{*j} = \text{norm}(\text{States}_{k}^{j}), \\ &\text{KL}^{j} = \sum_{c \in \mathcal{C}} \text{States}_{k}^{*j}(c) \log \frac{\text{States}_{k}^{*j}(c)}{\text{States}_{k/2}^{*j}(c)} + 1, \end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

where norm( $\cdot$ ) denotes the normalization operation.

Coarse-grained Changes However, focusing solely on fine-grained changes might not capture the 235 overall knowledge state changes of the learners. The Zone of Proximal Development (ZPD) is a 236 core concept in Dynamic Assessment Theory (Vygotsky & Cole, 1978). It refers to the gap between 237 a learner's current independent ability level and the potential level that could be reached with the 238 help of other mediums (e.g., ITS). It describes the overall changes in the learner's knowledge state 239 (*i.e.*, progress or regression). Inspired by this, we consider the overall correct rate at the intermediate 240 time-step k/2 as the lower limit of the ZPD, and the correct rate at the current time-step k as the 241 upper limit or near-upper limit of the ZPD. We use the rate of change as a quantitative indicator of 242 the ZPD<sup>j</sup> of learner  $s^{j}$ : 243

$$\operatorname{ZPD}^{j} = \left| \frac{\sum_{i=1}^{k} r_{i}^{j}}{\sum_{i=1}^{k} 1} - \frac{\sum_{i=1}^{k/2} r_{i}^{j}}{\sum_{i=1}^{k/2} 1} \right| \times \operatorname{len}_{j} \div \left( \frac{\sum_{i=1}^{k/2} r_{i}^{j}}{\sum_{i=1}^{k/2} 1} + 1 \right) + 1,$$
(4)

where  $len_j$  is the actual length of questions attempted by learner  $s^j$  ( $len_j \le k$ , when  $len_j < k$ , the missing sequence, *e.g.*, concepts sequence, is often padded with 0), which reflects the reliability of the ZPD<sup>j</sup>, with a larger  $len_j$  indicating more reliable results.

Finally, the controller assigns a value score to learner  $s^j$ :

core<sub>i</sub> = 
$$\mathrm{KL}^j \times \mathrm{ZPO}^j$$
. (5)

It can be observed that  $KL^{j}$  and  $ZPO^{j}$  are positive, which avoids any absolute impact on the score<sub>j</sub> when either one is 0. Notably,the controller can identify learners who have shown significant progress or regression, which is beneficial for teachers or ITS to pay further attention to them.

SC

## 3.2.2 GENERATOR

The generator can generate personalized dynamic parameters for learners determined by the controller based on real-time samples from different stages or groups, aiming to improve the adaptive generalization for continuously changing distributions. We first introduce the generator's feature extraction, then propose our designed state-adaptive attention, and finally discuss generating parameters through low-rank decomposition. For convenience, we have omitted the superscript of the learner numbers.

Feature Extraction. At time-step k, the generator takes  $\{(c_i, r_i)\}_{i=1}^k$  as input, considering the relative relationship between question difficulty and learner ability (Rasch, 1993) and inspired by the dual-tower model in recommendation systems (Huang et al., 2013; Covington et al., 2016), embedding the questions  $c_{1:k}$  and responses  $r_{1:k}$  into vector spaces  $Q_{1:k} \in \mathbb{R}^d$  and  $R_{1:k} \in \mathbb{R}^d$ , respectively (d is the dimension of the embedding). After non-linearization, features  $H_k^q \in \mathbb{R}^{d_{in}}$ and  $H_k^r \in \mathbb{R}^{d_{in}}$  ( $d_{in}$  is the input dimension of the dynamic layer) are extracted through a sequential feature extractor (SFE) (e.g., GRU): 

272  
273
$$\begin{cases}
H_k^q = SFE(Tanh(Q_{1:k}W_1 + b_1)), \\
H_k^r = SFE(Tanh(R_{1:k}W_2 + b_2)),
\end{cases}$$
(6)

where  $W_1 \in \mathbb{R}^{d \times d_{in}}$ ,  $W_2 \in \mathbb{R}^{d \times d_{in}}$ ,  $b_1 \in \mathbb{R}^{d_{in}}$ ,  $b_2 \in \mathbb{R}^{d_{in}}$  are learnable parameters in the projection-tion layer. Tanh( $\cdot$ ) is the activation function. 

State-adaptive Attention (SAA). SAA is responsible for adaptive generalization of the extracted question and response features, considering both change in concept correct rate (*i.e.*, difficulty) and the time of the change in knowledge state. Intuitively, the greater the change in difficulty, indicating more significant progress or regression, and the longer the time since the last response, the more likely a sudden change in knowledge state can occur. Such positions should receive more attention. Therefore, the definition of SAA is as follows: 

$$\begin{cases} SAA(X_k) = Concat(head_1, \cdots, head_h)W_h, \\ head_i = Attention^*(Q = X_k^{/h}, K = X_k^{/h}, V = X_k^{/h}), \\ Attention^*(Q, K, V) = softmax^*(X = \frac{QK^T}{\sqrt{d/h}})V, \\ softmax^*(X) = attn_w(c_{1:k}, r_{1:k}, t_{1:k}) \cdot softmax(X), \end{cases}$$
(5)

$$\mathsf{Lattn}_w(c_{1:k}, r_{1:k}, t_{1:k}) = \mathsf{dist}_d(c_{1:k}, r_{1:k}) \cdot \mathsf{dist}_t(c_{1:k}, t_{1:k}),$$

where h is the number of attention heads and  $W_h \in \mathbb{R}^{d_{in} \times d_{in}}$ .  $X_k^{/h}$  represents splitting the  $d_{in}$  di-mensions of  $X_k$  into h parts. dist<sub>d</sub> and dist<sub>t</sub> represent the changes in difficulty and time, respectively. At position  $i \in [1, k]$ , dist<sub>d</sub> $(c_i, r_i)$  and dist<sub>t</sub> $(c_i, t_i)$  are respectively: 

$$\begin{cases} 1, & \text{if } i = 1, \\ \left(\frac{\sum_{j=1}^{i} r_j [c_j = c_i]}{\sum_{j=1}^{i} 1 [c_j = c_i]} - \frac{\sum_{j=1}^{i-1} r_j [c_j = c_i]}{\sum_{j=1}^{i-1} 1 [c_j = c_i]}\right) + 1. & \text{else} \\ 1, & \text{if } j = \max\{k \mid k < i \text{ and } c_k = c_i\} = \emptyset, \\ \frac{t_i - t_j}{t_i - t_1}. & \text{else} \end{cases}$$
(8)

Finally, the representations  $S_k^q$  and  $S_k^r$  of question difficulty and learner ability are obtained by SAA:  $S_k^q = SAA(H_k^q), S_k^r = SAA(H_k^r),$ (9)

where  $S_{k}^{q}$  and  $S_{k}^{r}$  characterize the difficulty distribution of questions and the ability distribution of learners, respectively, based on real-time data from the current stage or group. SAA is the core component of the generator, and we will further discuss its importance in Sec. 4.5. 

Low-rank Decomposition. Before performing low-rank decomposition on the parameters, the learned question difficulty  $S_k^q$  and learner ability  $S_k^r$  are uniformly expressed as the generalized information feature  $S_k$  that characterizes the interaction distribution of learners: 

$$S_k = S_k^q + S_k^r,\tag{10}$$

Finally, parameters (*i.e.*, weight and bias) are generated through  $S_k$  for the dynamic layer:

$$\begin{cases} \text{weight} = S_k W_w + b_w, \\ \text{bias} = S_k W_b + b_b, \end{cases}$$
(11)

where  $W_w \in \mathbb{R}^{d_{in} \times (d_{in} \times d_{out})}$ ,  $b_w \in \mathbb{R}^{d_{in} \times d_{out}}$ ,  $W_b \in \mathbb{R}^{d_{in} \times d_{out}}$ ,  $b_b \in \mathbb{R}^{d_{out}}$  are learnable pa-rameters.  $d_{out}$  is the output dimension of the dynamic layer. However, it can be observed that the parameter size of  $W_w$  is too large, which increases computational resources and the risk of overfit-ting. Inspired by LoRA (Hu et al., 2021),  $W_w$  is decomposed into low-rank matrices to obtain the final weight: 

$$weight = S_k W_{w_1} W_{w_2} + b_w, \tag{12}$$

where  $W_{w_1} \in \mathbb{R}^{d_{in} \times \text{rank}}$ ,  $W_{w_2} \in \mathbb{R}^{\text{rank} \times (d_{in} \times d_{out})}$  are learnable parameters, and rank  $\ll d_{in}$  is a very small value (e.g., 1). In Sec. 4.5, we will further analyze the effects of different rank. 

It's noted that the generator can generate parameters for the dynamic layer, given the input dimension  $d_{in}$  and output dimension  $d_{out}$ . In our experiments, the dynamic layer defaults to the output layer of the KT model.

#### 324 3.2.3 MODEL TRAINING 325

326

331

332

333

334

335

336

337 338

339 340

341

349 350

361

All learnable parameters are trained by minimizing the binary cross-entropy between  $r_i$  and  $\hat{r}_i$ , *i.e.*,

$$\mathcal{L} = -\sum_{i=1}^{k} r_i \log(\hat{r}_i) + (1 - r_i) \log(1 - \hat{r}_i).$$
(13)

#### EXPERIMENTS 4

In this section, we demonstrate the superiority of our proposed Cuff-KT and the impact of its different components through experiments. Specifically, the experimental evaluation is divided into (i) the controllability of parameter generation (Sec. 4.2), (ii) prediction accuracy, quantifying the effectiveness of tackling RLPA (Sec. 4.3), (iii) the application of Cuff-KT (Sec. 4.4), and (iv) the impact of dual-tower modeling, SFE, SAA, and low-rank decomposition in Cuff-KT (Sec. 4.5).

# 4.1 EXPERIMENTAL SETUP

# 4.1.1 DATASETS

342 We conduct extensive experiments on a classic dataset (assist15 (Feng et al., 2009)) and two recently proposed benchmark datasets (comp (Hu et al., 2023) and xes3g5m (Liu et al., 2024)). Following 343 the data preprocessing method outlined in (Lee et al., 2022b), we exclude learners with fewer than 344 five interactions and all interactions involving nameless concepts. Since a question may involve 345 multiple concepts, we convert the unique combinations of concepts within a single question into 346 a new concept. Table 1 provides a statistical overview of these datasets. It's noted that the large 347 datasets (comp and xes3g5m) are randomly sampled 5000 learners. 348

See Appendix A.1 for a detailed description of the datasets.

#### 4.1.2 BASELINES 351

352 We select a classic KT model 353 (DKT (Piech et al., 2015)) and 354 two recently proposed state-of-the-355 art models (AT-DKT (Liu et al., 356 2023a) and DIMKT (Shen et al., 357 2022)) as the backbone models for

Table 1: Statistics of 3 datasets.

Datasets	#learners	#questions	#concepts	#interactions
assist15	17,115	100	100	676,288
comp	5,000	7,460	445	668,927
xes3g5m	5,000	7,242	1,221	1,771,657

358 optimization. We compare Cuff-KT with these three backbone models and three classic fine-tuning based methods: Full Fine-tuning (FFT), Adapter-based tuning (Adapter) (Houlsby et al., 2019), and 359 Bias-term Fine-tuning (BitFit) (Zaken et al., 2021). 360

See detailed introductions to the backbone models and baselines in Appendix A.3. 362

# 4.1.3 IMPLEMENTATION

364 We implement all models using Pytorch on a Linux server with two GeForce RTX 3090s. We 365 used the Adam optimizer with a learning rate of 0.001, and a batch size of 512. The embedding 366 dimension for all models is fixed at 32. The rank of the generator in Cuff-KT is set to 1. We split 367 the historical interactions into training, validation, and test sets (7:2:1) based on timestamps and 368 groups, respectively. An early stopping strategy is applied if the AUC on the validation set does 369 not increase for 10 epochs. The experiments are repeated 5 times under random seeds 0 to 4 and 370 the average performance is reported. Following the previous works (Piech et al., 2015; Shen et al., 371 2022; Liu et al., 2023a), the evaluation metrics include Area Under the ROC Curve (AUC) and Root 372 Mean Square Error (RMSE).

373

#### 374 4.2 CONTROLLABLE PARAMETER GENERATION 375

According to (Lv et al., 2023a), anomaly detection algorithms can be used to detect distribution 376 changes over time. We select four representative anomaly detection algorithms from pyod li-377 brary (Zhao et al., 2019) as comparison baselines for the controller in Cuff-KT: LOF (Breunig et al.,



Figure 4: Performance comparison of Cuff-KT and anomaly detection algorithms at different frequencies.

Table 2: Performance comparison between different methods under intra-learner shift. The **best** result is in bold and the <u>next best</u> is underlined. \* and \*\* indicate that the improvements over the strongest baseline are statistically significant, with p < 0.05 and p < 0.01, respectively.

$Dataset \rightarrow$		assist15	assist15					xes3g5m		
$\overline{Method {\downarrow} \backslash Metric} {\rightarrow}$	AUC ↑	$\text{RMSE}\downarrow$	Time Cost $\downarrow$	AUC	RMSE	Time Cost	AUC	RMSE	Time Cost	
DKT	0.7058	0.4107	0ms	0.6990	0.3613	0ms	0.6633	0.4129	0ms	
+FFT	<u>0.7063</u>	0.4071	≥17,200ms	<u>0.7066</u>	0.3594	≥18,300ms	<u>0.7116</u>	0.3992	≥33,600ms	
+Adapter	0.6749	0.4242	≥16,600ms	0.6634	0.3714	$\geq$ 17,700ms	0.6467	0.4275	$\geq$ 36,100ms	
+BitFit	0.7054	0.4080	≥16,300ms	0.7039	0.3599	≥14,100ms	0.6841	0.4105	≥32,600ms	
+Cuff-KT	0.8130**	0.3773**	$\geq$ 419ms	0.7834**	0.3459**	$\geq$ 435ms	0.7176	0.3931	$\geq$ 1,211ms	
AT-DKT	0.6981	0.4106	0ms	0.6922	0.3621	0ms	0.6437	0.4228	0ms	
+FFT	<u>0.7005</u>	0.4083	≥126,400ms	0.7020	0.3602	≥95,000ms	0.6918	<u>0.4068</u>	≥176,000ms	
+Adapter	0.6588	0.4287	$\geq$ 125,100ms	0.6443	0.3878	$\geq$ 88,800ms	0.6276	0.4351	$\geq$ 168,300ms	
+BitFit	0.6989	0.4094	$\geq$ 121,300ms	0.6990	0.3608	≥91,300ms	0.6668	0.4178	≥169,300ms	
+Cuff-KT	0.8335**	0.3714**	$\geq$ 236ms	0.7869**	0.3435**	$\geq$ 254ms	0.7133**	0.4009*	$\geq$ 784ms	
DIMKT	0.7055	0.4080	0ms	0.7934	0.3404	0ms	0.8322	0.3402	0ms	
+FFT	0.7072	0.4063	≥270,900ms	0.8000	0.3375	≥205,200ms	0.8366	0.3383	≥377,800ms	
+Adapter	0.6507	0.4387	≥410,000ms	0.7526	0.3671	≥278,500ms	0.7929	0.3696	≥509,200ms	
+BitFit	0.7082	0.4061	$\geq$ 263,500ms	0.7972	0.3382	≥199,400ms	0.8369	0.3381	$\geq$ 347,800ms	
+Cuff-KT	0.8322**	0.3710*	$\geq$ 232ms	0.8380**	0.3297**	$\geq 175 ms$	0.8540*	0.3347*	$\geq$ 239ms	

2000), PCA (Shyu et al., 2003), IForest (Liu et al., 2008), and ECOD (Li et al., 2022), and use
AUC as the evaluation metric. Detailed descriptions of these four algorithms can be found in the
Appendix A.2. Figure 4 shows the performance results under intra-learner shift when the controller
selects learners with different frequencies for the generator.

We can see that anomaly detection algorithms (especially IForest and ECOD) consistently outperform the random selection, demonstrating the correctness of using anomaly detection algorithms to
detect distribution changes. Moreover, our Cuff-KT generally performs better than these algorithms,
indicating that Cuff-KT is more capable of identifying learners whose model generalization deteriorates due to distribution changes. We attribute Cuff-KT's breakthrough to the Dynamic Assessment
Theory (Vygotsky & Cole, 1978), which we further analyze in the Appendix A.4.

### 4.3 TUNING-FREE AND FAST PREDICTION

Under this setting, the generator in Cuff-KT generates parameters for all learners independently of
the controller. In our setup, we attempt to divide learners into different groups based on the degree
of change in their knowledge states. We use DKT to encode each learner's interaction history and
choose the distance (*e.g.*, KL divergence) between the prediction distributions for each concept at the

434	strongest baseling	ne are stat	istically	significant,	with $p <$	<0.05 an	d <i>p</i> <0.01,	respecti	vely.	
435	$Dataset \rightarrow$		assist15	5		comp			xes3g5m	
436	$Method{\downarrow}\backslash Metric{\rightarrow}$	AUC ↑	$RMSE \downarrow$	Time Cost $\downarrow$	AUC	RMSE	Time Cost	AUC	RMSE	Time Cost
437	DKT	0.7075	0.4363	0ms	0.6681	0.4355	0ms	0.7907	0.4329	0ms
438	+FFT	0.7137	$\frac{0.4339}{0.4456}$	$\geq 18,800$ ms	0.6839 0.6461	$\frac{0.4310}{0.4438}$	$\geq$ 3,600ms	0.7990 0.7646	$\frac{0.4166}{0.4427}$	$\geq$ 4,400ms
439	+BitFit	0.7119	0.4349	$\geq 17,000$ ms $\geq 17,200$ ms	0.6734	0.44326	$\geq$ 3,200ms $\geq$ 3,100ms	0.7905	0.4323	$\geq$ 4,900ms
440	+Cuff-KT	0.7365*	0.4302	$\geq$ 355ms	0.6937**	0.4294*	$\geq$ 96ms	0.8004	0.4158	$\geq 123 ms$
441	AT-DKT	0.7030	0.4389	0ms	0.6587	0.4375	Oms	0.7868	0.4370	0ms
442	+FF1 +Adapter	$\frac{0.7104}{0.6708}$	$\frac{0.4355}{0.4520}$	$\geq$ 74,700ms $\geq$ 55 300ms	$\frac{0.6751}{0.6253}$	$\frac{0.4312}{0.4498}$	$\geq 20,300$ ms $\geq 18,100$ ms	0.7916 0.7643	$\frac{0.4242}{0.4457}$	$\geq 21,300$ ms $\geq 23,500$ ms
443	+BitFit	0.7076	0.4367	$\geq$ 59,100ms	0.6666	0.4334	$\geq$ 18,600ms	0.7860	0.4352	$\geq$ 19,300ms
444	+Cuff-KT	0.7348**	0.4316*	$\geq$ 170ms	0.6919**	0.4303*	≥64ms	0.7959*	0.4183*	≥110ms
115	DIMKT	0.7134	0.4350	0 ms >173 500 ms	0.7556	0.4118	0 ms >52 500 ms	0.8255	0.4088	0ms
445	+Adapter	0.6648	$\frac{0.4320}{0.4577}$	$\geq 173,300 \text{ms}$ $\geq 217,100 \text{ms}$	0.7017	$\frac{0.4097}{0.4465}$	$\geq$ 32,500ms $\geq$ 82,500ms	0.7467	0.4618	$\geq$ 43,200ms $\geq$ 81,800ms
440	+BitFit	0.7144	0.4334	$\geq$ 154,400ms	0.7563	0.4110	≥50,800ms	0.8254	0.4084	≥48,600ms
447	+Cutt-KT	0.7425**	0.4296	$\geq$ 203ms	0.7657**	0.4057*	$\geq$ 64ms	0.8309	0.4009	$\geq /2ms$
448	0.86 - assi	st15	0.40	0.87	comp	,	0.36 0.88 -	x	es3g5m	0.42
449	Cuff-KT + FFT Cuff-KT			Cuff-KT + FFT Cuff-KT		0.861		Cuff-KT + FFT Cuff-KT	0.865	, 4
450	0.84 RMSE	0.846		0.85 RMSE			0.84	RMSE	0.401	- 0.40
451	0.381 0.834	0.832	0.38	0.346	0.344	0.838	0.34	0.393	0.387	
452	y 0.82 •	0.371 0.	371	0.83	0.336		5 C 2			- 0.38
453	4		2	4 0.81 AUC	0.907	0.330	0.77 A	uc	58	- 0.36
454	AUC 0.80 - 0.803	0.3	60 0.36	0.804	0.007	-	0.32	152		
455	0.793	0.352		0.79	0.787	0.314	0.74	719		- 0.34 0.335 0.332
456	0.78		0.34	0.77			0.30 0.70	0.7	13	0.32
457	DKI AT- Moi	dels	<b>.</b>		Models				Models	UMNI
458		]	Figure 5	: Cutt-KT+	FFT und	er intra-l	earner shif	t.		

Table 3: Performance comparison between different methods under inter-learner shift. The **best** result is in bold and the <u>next best</u> is underlined. \* and \*\* indicate that the improvements over the strongest baseline are statistically significant, with p < 0.05 and p < 0.01, respectively.

intermediate and current timestamps as the basis for division. Tables 2 and 3 show the performance comparison between different methods under intra-learner shift and inter-learner shift. Overall, our Cuff-KT effectively tackles the RLPA task with significant advantages. We can observe:

- Overall, compared to baseline methods, our Cuff-KT generally performs best on all metrics across all datasets. This performance improvement can be attributed to Cuff-KT's parameter generation approach, which dynamically updates the model to capture distribution dynamics rather than statically considering interactions in the test data, enhancing the KT model's dynamic adaptability.
- Compared to the backbone, the time cost caused by Cuff-KT is significantly smaller than fine-tuning-based methods. This is because Cuff-KT updates model parameters only through feedforward computation, without the need for a retraining process.
- Adapter fine-tuning performs poorly and even leads to performance degradation, as it is heavily affected by task complexity and model scale (He et al., 2021; Karimi Mahabadi et al., 2021), ultimately resulting in overfitting.
- Although FFT and BitFit fine-tuning methods generally improve the performance of the backbone, especially FFT based on DKT showing a 0.483 increase in AUC metric on the xes3g5m dataset under intra-learner shift, the time cost caused is non-negligible in real-world scenarios.
- 477 478 479

480

459

460

461

462 463

464

465

466

467

468

469

470

471

472

473

474

475

476

432

433

# 4.4 FLEXIBLE APPLICATION

Thanks to the independence of the generator in our Cuff-KT from fine-tuning based methods, we attempt to combine Cuff-KT with FFT. The results in terms of AUC and RMSE under intra-learner shift and inter-learner shift are shown in Figure 5 and Figure 8 in the Appendix A.4, respectively. As can be seen from the figures, on different backbone models and across all datasets, the performance still shows a significant improvement after combining Cuff-KT with FFT. This is because FFT can learn different distributions from the recent data, facilitating Cuff-KT's smooth transition to the



Figure 6: Performance on AUC for different rank under intra-learner shift.

distribution in the test data. This combination provides a reference for flexibly fine-tuning models in special real-world scenarios where real-time requirements are not high.

Moreover, the generator in Cuff-KT can flexibly generate parameters for or insert into any layer of the KT model. This inspires us to consider how the generator can choose the position and network structure for generation or insertion. Due to space limitations, we leave this as a direction for future research.

502 504

505

496

497

498 499

500

501

#### 4.5 ABLATION STUDY

We systematically examine the im-506 pact of key components in Cuff-KT 507 based on DKT by constructing four 508 variants under intra-learner shift. 509 "w/o. Dual" indicates that question 510 and response embeddings are fused 511 (e.g., by summation) after embed-512 ding. "w/o. SFE" means the SFE 513 component is omitted, "w/o. SAA" 514 means omitting the SAA compo-

Table 4: The performance of different variants in Cuff-KT.

$\text{Dataset} \rightarrow$	ass	ist15	co	mp	xes3g5m		
$Metric \rightarrow$	AUC↑	RMSE↓	AUC	RMSE	AUC	RMSE	
Cuff-KT	0.8130	0.3773	0.7834	0.3459	0.7176	0.3931	
w/o. Dual	0.7013	0.4126	0.7245	0.3693	0.7088	0.4094	
w/o. SFE	0.7706	0.3925	0.7204	0.3612	0.6896	0.4140	
w/o. SAA	0.7000	0.4141	0.6877	0.3640	0.6716	0.4212	
w. SHA	0.7810	0.3844	0.6924	0.3629	0.6767	0.4185	

515 nent, and "w. SHA" means SAA is replaced with standard multi-head attention. From Table 4, 516 we can easily observe: (1) Cuff-KT outperforms all variants, especially when the SAA component 517 is removed, the predictive performance generally decreases the most, while Cuff-KT with standard multi-head attention comes next, empirically validating that our designed SAA component can ef-518 fectively achieve adaptive generalization. (2) Cuff-KT's performance is very low when the SFE 519 component is removed or dual modeling is not employed. We believe this is because Cuff-KT can 520 successfully extract question features and learner response features and effectively learn the dif-521 ficulty distribution of current questions and the ability distribution of learners based on real-time 522 data. 523

Additionally, we study the effects of different rank under intra-learner shift. The performance on 524 AUC of different rank under intra-learner shift and the parameter size of the generator in Cuff-KT 525 are shown in Figure 6 and Table 5 in the Appendix A.4, respectively. In Figure 6, after low-rank 526 decomposition (rank  $\neq 0$ ), the performance on AUC generally improves, and the effects of different 527 rank are inconsistent across different datasets. In Table 5, the parameter size of the generator in-528 creases with the rank, indicating that by adjusting different ranks, an effective balance between the 529 performance and resource consumption of Cuff-KT can be achieved. 530

531 532

533

#### CONCLUSION 5

534 Our paper aims to tackle the RLPA task in KT by proposing a controllable, tuning-free, fast, and flexible method called Cuff-KT to improve adaptability of KT models in real-world scenarios. We 536 decompose the RLPA task to be solved into two sub-issues: intra-learner shift and inter-learner shift, 537 and design a parameter generator capable of generate personalized parameters based on the current stage or group, thereby achieving adaptive generalization. In instance validations across multiple 538 KT models, Cuff-KT exhibits superior performance in adapting to rapidly changing distributions, avoiding the overfitting and high time cost challenges inherent in fine-tuning based methods.

540	REFERENCES
541	THE ENERGES

- Ghodai Abdelrahman, Qing Wang, and Bernardo Nunes. Knowledge tracing: A survey. ACM
   *Computing Surveys*, 55(11):1–37, 2023.
- Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying densitybased local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pp. 93–104, 2000.
- Albert T Corbett and John R Anderson. Knowledge tracing: Modeling the acquisition of procedural
   knowledge. *User modeling and user-adapted interaction*, 4:253–278, 1994.
- Paul Covington, Jay Adams, and Emre Sargin. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*, pp. 191–198, 2016.
- Mingyu Feng, Neil Heffernan, and Kenneth Koedinger. Addressing the assessment challenge with
   an online system that tutors as it assesses. *User modeling and user-adapted interaction*, 19:243–266, 2009.
- Aritra Ghosh, Neil Heffernan, and Andrew S. Lan. Context-aware attentive knowledge tracing. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20, pp. 2330–2339, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3403282. URL https://doi. org/10.1145/3394486.3403282.
- Xiaopeng Guo, Zhijie Huang, Jie Gao, Mingyu Shang, Maojing Shu, and Jun Sun. Enhancing knowledge tracing via adversarial training. In *Proceedings of the 29th ACM International Conference on Multimedia*, MM '21, pp. 367–375, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450386517. doi: 10.1145/3474085.3475554. URL https://doi.org/10.1145/3474085.3475554.
- Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a
   unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*, 2021.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
  and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Liya Hu, Zhiang Dong, Jingyuan Chen, Guifeng Wang, Zhihua Wang, Zhou Zhao, and Fei Wu. Ptadisc: A cross-course dataset supporting personalized learning in cold-start scenarios. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. Learning deep
   structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, pp. 2333–2338, 2013.
- Firuz Kamalov, David Santandreu Calonge, and Ikhlaas Gurrib. New era of artificial intelligence in education: Towards a sustainable multifaceted revolution. *Sustainability*, 15(16), 2023. ISSN 2071-1050. doi: 10.3390/su151612451. URL https://www.mdpi.com/2071-1050/15/ 16/12451.
- Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank
   hypercomplex adapter layers. Advances in Neural Information Processing Systems, 34:1022–1035, 2021.

608

- 594 Wonsung Lee, Jaeyoon Chun, Youngmin Lee, Kyoungsoo Park, and Sungrae Park. Contrastive 595 learning for knowledge tracing. In Proceedings of the ACM Web Conference 2022, WWW '22, 596 pp. 2330–2338, New York, NY, USA, 2022a. Association for Computing Machinery. ISBN 597 9781450390965. doi: 10.1145/3485447.3512105. URL https://doi.org/10.1145/ 598 3485447.3512105.
- Wonsung Lee, Jaeyoon Chun, Youngmin Lee, Kyoungsoo Park, and Sungrae Park. Contrastive 600 learning for knowledge tracing. In Proceedings of the ACM Web Conference 2022, pp. 2330– 2338, 2022b. 602
- Zheng Li, Yue Zhao, Xiyang Hu, Nicola Botta, Cezar Ionescu, and George H Chen. Ecod: Unsu-603 pervised outlier detection using empirical cumulative distribution functions. IEEE Transactions 604 on Knowledge and Data Engineering, 35(12):12181-12193, 2022. 605
- 606 Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In 2008 Eighth IEEE Interna-607 tional Conference on Data Mining, pp. 413–422, 2008. doi: 10.1109/ICDM.2008.17.
- Qi Liu, Shuanghong Shen, Zhenya Huang, Enhong Chen, and Yonghe Zheng. A survey of knowl-609 edge tracing. arXiv preprint arXiv:2105.15106, 2021. 610
- 611 Zitao Liu, Qiongqiong Liu, Jiahao Chen, Shuyan Huang, Boyu Gao, Weiqi Luo, and Jian Weng. 612 Enhancing deep knowledge tracing with auxiliary tasks. In Proceedings of the ACM Web Con-613 ference 2023, WWW '23, pp. 4178–4187, New York, NY, USA, 2023a. Association for Computing Machinery. ISBN 9781450394161. doi: 10.1145/3543507.3583866. URL https: 614 //doi.org/10.1145/3543507.3583866. 615
- 616 Zitao Liu, Qiongqiong Liu, Jiahao Chen, Shuyan Huang, and Weiqi Luo. simplekt: A simple but 617 tough-to-beat baseline for knowledge tracing, 2023b.
- Zitao Liu, Qiongqiong Liu, Teng Guo, Jiahao Chen, Shuyan Huang, Xiangyu Zhao, Jiliang Tang, 619 Weiqi Luo, and Jian Weng. Xes3g5m: A knowledge tracing benchmark dataset with auxiliary 620 information. Advances in Neural Information Processing Systems, 36, 2024. 621
- 622 Zheqi Lv, Zhengyu Chen, Shengyu Zhang, Kun Kuang, Wenqiao Zhang, Mengze Li, Beng Chin 623 Ooi, and Fei Wu. Ideal: Toward high-efficiency device-cloud collaborative and dynamic recommendation system. arXiv preprint arXiv:2302.07335, 2023a. 624
- 625 Zheqi Lv, Wenqiao Zhang, Shengyu Zhang, Kun Kuang, Feng Wang, Yongwei Wang, Zhengyu 626 Chen, Tao Shen, Hongxia Yang, Beng Chin Ooi, et al. Duet: A tuning-free device-cloud collabo-627 rative parameters generation framework for efficient device model generalization. In Proceedings 628 of the ACM Web Conference 2023, pp. 3077-3085, 2023b. 629
- Hiromi Nakagawa, Yusuke Iwasawa, and Yutaka Matsuo. Graph-based knowledge tracing: model-630 ing student proficiency using graph neural network. In IEEE/WIC/ACM International Conference 631 on Web Intelligence, pp. 156-163, 2019. 632
- 633 Shalini Pandey and George Karypis. A self-attentive model for knowledge tracing. arXiv preprint 634 arXiv:1907.06837, 2019.
- 635 Shalini Pandey and Jaideep Srivastava. Rkt: Relation-aware self-attention for knowledge trac-636 ing. In Proceedings of the 29th ACM International Conference on Information & Knowledge 637 Management, CIKM '20, pp. 1205-1214, New York, NY, USA, 2020. Association for Com-638 puting Machinery. ISBN 9781450368599. doi: 10.1145/3340531.3411994. URL https: 639 //doi.org/10.1145/3340531.3411994. 640
- Chris Piech, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, 641 Leonidas J Guibas, and Jascha Sohl-Dickstein. Deep knowledge tracing. In C. Cortes, 642 N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), Advances in Neural In-643 formation Processing Systems, volume 28. Curran Associates, Inc., 2015. URL 644 https://proceedings.neurips.cc/paper\_files/paper/2015/file/ 645 bac9162b47c56fc8a4d2a519803d51b3-Paper.pdf. 646
- Sidney L Pressey. A simple apparatus which gives tests and scores-and teaches. Sch. & Soc., 23: 647 373-376, 1926.
  - 12

648 649	Georg Rasch. Probabilistic models for some intelligence and attainment tests. ERIC, 1993.
650	Shuanghong Shen Zhenya Huang Oi Liu Yu Su Shijin Wang and Enhong Chen Assessing
651	student's dynamic knowledge state by exploring the question difficulty effect. In <i>Proceedings</i>
652	of the 45th International ACM SIGIR Conference on Research and Development in Informa-
653	tion Retrieval, SIGIR '22, pp. 427-437, New York, NY, USA, 2022. Association for Com-
654	puting Machinery. ISBN 9781450387323. doi: 10.1145/3477495.3531939. URL https:
655	//doi.org/10.1145/3477495.3531939.
656	Mai Ling Chuy, Chu, Ching Chan, Kanalani Saninganakam, and LiWy, Chang. A naval anomaly
657	detection scheme based on principal component classifier. In <i>Proceedings of the IEEE foundations</i>
658	and new directions of data mining workshop np 172–179 IFFE Press Piscataway NI USA
659	2003.
660	
661	Lev Semenovich Vygotsky and Michael Cole. <i>Mind in society: Development of higher psychologi</i> -
662	cal processes. Harvard university press, 1978.
663	Cheryl Sze Yin Wong and Savitha Ramasamy. Architectural adaptation and regularization of atten-
664	tion networks for incremental knowledge tracing. In <i>Proceedings of the 14th Learning Analytics</i>
665	and Knowledge Conference, pp. 757–762, 2024.
666	
667	Cheryl Sze Yin Wong, Guo Yang, Nancy F Chen, and Ramasamy Savitha. Incremental context
668	aware allentive knowledge tracing. In ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 2003–2007, IEEE 2022
669	Acoustics, speech and signal Processing (TCASSP), pp. 5995-5997. IEEE, 2022.
670	Zhengyi Yang, Xiangnan He, Jizhi Zhang, Jiancan Wu, Xin Xin, Jiawei Chen, and Xiang Wang. A
671	generic learning framework for sequential recommendation with distribution shifts. In Proceed-
672	ings of the 46th International ACM SIGIR Conference on Research and Development in Informa-
673	tion Retrieval, pp. 331–340, 2023.
674	Chun-Kit Yeung and Dit-Yan Yeung. Addressing two problems in deep knowledge tracing via
675	prediction-consistent regularization. In Proceedings of the Fifth Annual ACM Conference on
676	Learning at Scale, L@S '18, New York, NY, USA, 2018. Association for Computing Machin-
677	ery. ISBN 9781450358866. doi: 10.1145/3231644.3231647. URL https://doi.org/10.
678	1145/3231644.3231647.
679	Elad Ben Zaken, Shauli Rayfogel, and Yoay Goldberg, Bitfit: Simple parameter-efficient fine-tuning
680	for transformer-based masked language-models. arXiv preprint arXiv:2106.10199, 2021.
681	
682	Jiani Zhang, Xingjian Shi, Irwin King, and Dit-Yan Yeung. Dynamic key-value memory networks
683	for knowledge tracing. In Proceedings of the 20th international conference on world wide web,
684	pp. 703–774, 2017.
685	Yue Zhao, Zain Nasrullah, and Zheng Li. Pyod: A python toolbox for scalable outlier detection.
686	Journal of machine learning research, 20(96):1–7, 2019.
687	Hansi Zhan Dahart Damlar Chadan M We and Alam Thing Chattan David at a 111 - 1
688	Hanqi Zhou, Robert Bamier, Charley M Wu, and Alvaro Tejero-Cantero. Predictive, scalable and intermetable knowledge tracing on structured domains. arXiv preprint arXiv:2402.12170.2024
689	interpretable knowledge tracing on structured domains. arXiv preprint arXiv:2403.15179, 2024.
690	
691	
692	
693	
694	
695	
696	
697	
698	
599	
700	
/ U I	

# 702 A APPENDIX

704

705

706

711 712

713

714 715

716

717

718

719 720

721

722

723 724

726

727 728

729

730

731

732

733

734

735

737

738

739

740

741

742

In this document, we include the description of the datasets (A.1), an overview to anomaly detection algorithms (A.2), an introduction to the backbone models and baselines (A.3), and additional experimental results (A.4), which we are unable to include in the main paper due to space limitations.

- 708 A.1 DATASETS 709
- <sup>710</sup> Here, we describe the datasets (assist15, comp, xes3g5m) used for evaluation.
  - **assist15**<sup>3</sup> (Feng et al., 2009): The assist15 dataset, is collected from the ASSISTments platform in the year of 2015. It includes a total of 708,631 interactions involving 100 distinct concepts from 19,917 learners.
  - **comp**<sup>4</sup> (Hu et al., 2023): The comp dataset, is part of the PTADisc, which encompasses a wide range of courses from the PTA platform. PTADisc includes data from 74 courses, involving 1,530,100 learners and featuring 4,504 concepts, 225,615 questions, as well as an extensive log of over 680 million learner responses. The comp dataset is specifically selected for KT task in Computational Thinking course.
  - **xes3g5m**<sup>5</sup> (Liu et al., 2024): The xes3g5m dataset incorporates rich auxiliary information such as tree-structured concept relationships, question types, textual contents, and learner response timestamps and includes 7,652 questions and 865 concepts, with a total of 5,549,635 interactions from 18,066 learners.
- 725 A.2 ANOMALY DETECTION ALGORITHMS

We compare the controller of Cuff-KT with the following anomaly detection algorithms:

- LOF (Breunig et al., 2000): LOF quantifies the local outlier degree of samples by calculating a score. This score reflects the ratio of the average density of the local neighborhood around a sample point to the density at the location of that sample point. A ratio significantly greater than 1 indicates that the density at the sample point's location is much lower than the average density of its surrounding neighborhood, suggesting that the point is more likely to be a local outlier.
- PCA (Shyu et al., 2003): After performing eigenvalue decomposition, the eigenvectors obtained from PCA reflect different directions of variance change in network traffic data, while eigenvalues represent the magnitude of variance in the corresponding directions. Thus, the eigenvector associated with the largest eigenvalue represents the direction of maximum variance in network traffic data, while the eigenvector associated with the smallest eigenvalue represents the direction of minimum variance. If an individual network connection sample exhibits characteristics inconsistent with the overall network traffic sample, such as deviating significantly from other normal connection samples in certain directions, it may indicate that this connection sample is an outlier.
- 743 • IForest (Liu et al., 2008): IForest employs an innovative anomaly isolation method to 744 identify anomalous samples by constructing a binary tree structure (called an Isolation Tree or iTree). Unlike traditional methods, IForest does not build a model of normal samples, but 745 instead directly isolates anomalies. In this process, anomalous samples tend to be isolated 746 more quickly and thus are positioned closer to the root node in the tree, while normal 747 samples are isolated deeper in the tree. By constructing multiple iTrees (typically T trees), 748 the average path length from anomalies to the root node is significantly shorter than that of 749 normal points, and this characteristic is used for anomaly detection. This approach excels 750 in handling large-scale datasets and high-dimensional data, with the advantages of linear 751 time complexity and low memory requirements.

<sup>752</sup> \_

<sup>&</sup>lt;sup>3</sup>https://sites.google.com/site/assistmentsdata/datasets/

<sup>754 2015-</sup>assistments-skill-builder-data

<sup>&</sup>lt;sup>4</sup>https://github.com/wahr0411/PTADisc

<sup>&</sup>lt;sup>5</sup>https://github.com/ai4ed/XES3G5M

756		• ECOD (Li et al., 2022): ECOD is a novel unsupervised anomaly detection algorithm. Its
757		core idea stems from the definition of outliers—typically rare events occurring in the tails
758		of a distribution. The algorithm cleverly uses empirical cumulative distribution functions
759		(ECDF) to estimate the joint cumulative distribution function of the data, thereby calculat-
760		ing the probability of outliers. The uniqueness of ECOD lies in its avoidance of the slow
761		convergence problem of joint ECDF in high-dimensional data. The algorithm calculates
762		the univariate ECDF for each dimension separately, then estimates the degree of anomaly
763		for multidimensional data points through an independence assumption. This is done by
764		multiplying the estimated tail probabilities of all dimensions.
765		
766	A.3	BACKBONE MODELS AND BASELINES
767	Wai	estantiate a classic backhone model and two recently proposed SOTA models
768	we n	istantiate a classic backbone model and two recently proposed SOTA models.
769		• DKT (Piech et al., 2015): DKT is a seminal model that leverages Recurrent Neural Net-
770		works (RNNs), specifically utilizing a single layer LSTM, to directly model learners' learn-
771		ing processes and predict their performance.
772		• AT-DKT (Liu et al., 2023a): AT-DKT augments the original deep knowledge tracing model
773		by embedding two auxiliary learning tasks: one for predicting concepts and another for
774		assessing individualized prior knowledge. This integration aims to sharpen the model's
775		predictive accuracy and deepen its understanding of learner performance.
776		• <b>DIMKT</b> (Shen et al., 2022): DIMKT is designed to enhance the assessment of learners'
777		knowledge states by explicitly incorporating the difficulty level of questions and establishes
778		the relationship between learners' knowledge states and difficulty level during the practice
779		process.
780		
781	We c	ompare Cuff-KT with three classic fine-tuning based methods.
782		• Full Fine tuning (FFT): EET involves training all parameters of a model completely. It
783		usually has the highest potential for performance but it consumes the most resources takes
784		the longest time to train, and is prone to overfitting when the corpus is not large enough.
785		• Adapter-hased tuning (Adapter) (Houlsby et al. 2019): Adapter inserts downstream task
786		parameters, known as adapters, into each Transformer block of the pre-trained model. Each
787		adapter consists of two layers of MLP and an activation function, responsible for reducing
788		and increasing the dimensionality of the Transformer's representations. During fine-tuning,
789		the main model parameters are frozen, and only the task-specific parameters are trained.
790		Since the backbone models might not include a Transformer, in our experiments, it is re-
791		placed by linear layers.
792		• Bias-term Fine-tuning (BitFit) (Zaken et al., 2021): BitFit is a sparse fine-tuning method
793		that efficiently tunes only the parameters with bias, while all other parameters are fixed.
794		This method tends to be effective on small to medium datasets and can even compete with
795		other sparse fine-tuning methods on large datasets.
796		
797	A.4	Additional Experimental Results
798		
799		1. We further analyze the influence of different components of the controller in Cuff-KT under
800		nura-rearner snill, we instantiate DK1 on assist15, comp, and xes3g5m datasets. The AUC
801		when the controller removes ZPD ("w/o, ZPD" <i>i.e.</i> without considering coarse grained
802		changes in knowledge states) This indicates that considering coarse-grained knowledge
803		state changes is crucial, which aligns with reality, as in practical scenarios, a learner's
804		progress or regression is often judged by an overall score. Additionally, when ZPD does
805		not take into account actual length ("w/o. Rel.", i.e., without considering the reliability of
806		ZPD), the performance drops the second most. This is because when a learner has more
807		activity records, their knowledge state is more likely to experience drastic changes, and
808		such learners should receive more attention. On the other hand, when a learner has limited
809		records, the simulated changes in their knowledge state are less reliable and should be given

lower weight. When the controller does not consider fine-grained changes in the knowledge



- experience a major shift. However, such situations are relatively rare in reality.2. Figure 8 shows the performance of Cuff-KT combined with FFT under inter-learner shift
- in terms of AUC and RMSE.
- 3. Table 5 shows the parameter sizes (k) of the generator in Cuff-KT with different ranks under intra-learner shift.

Table 5: The parameter size (k) of the generator with different rank in Cuff-KT under intra-learner shift.

Dataset	Backhone	Rank						
Dataset		0	1	2	3	4		
	DKT	130.12	29.96	33.22	36.49	39.75		
assist15	AT-DKT	130.12	29.96	33.22	36.49	39.75		
	DIMKT	54.98	23.26	24.32	25.38	26.43		
	DKT	516.86	74.46	88.77	103.07	117.37		
comp	AT-DKT	516.86	74.46	88.77	103.07	117.37		
I I	DIMKT	66.02	34.30	35.36	36.42	37.47		
	DKT	1,386.76	174.57	213.70	252.84	291.97		
xes3g5m	AT-DKT	1,386.76	174.57	213.70	252.84	291.97		
2	DIMKT	90.85	59.14	60.19	61.25	62.30		