# A Survey of Knowledge Tracing

Qi Liu, Shuanghong Shen, Zhenya Huang, Enhong Chen, *Senior Member, IEEE* and Yonghe Zheng

**Abstract**—High-quality education is one of the keys to achieving a more sustainable world. The recent COVID-19 epidemic has triggered the outbreak of online education, which has enabled both students and teachers to learn and teach at home. Meanwhile, it is now possible to record and research a large amount of learning data using online learning platforms in order to offer better intelligent educational services. Knowledge Tracing (KT), which aims to monitor students' evolving knowledge state, is a fundamental and crucial task to support these intelligent services. Therefore, an increasing amount of research attention has been paid to this emerging area and considerable progress has been made. In this survey, we propose a new taxonomy of existing basic KT models from a technical perspective and provide a comprehensive overview of these models in a systematic manner. In addition, many variants of KT models have been proposed to capture more complete learning process. We then review these variants involved in three phases of the learning process: before, during, and after the student learning, respectively. Moreover, we present several typical applications of KT in different educational scenarios. Finally, we provide some potential directions for future research in this fast-growing field.

**Index Terms**—Knowledge Tracing; Intelligent Education; Educational Data Mining; Adaptive Learning; User Modeling

✦

## 1 INTRODUCTION

THROUGHOUT the world, education has been found to exhibit a close and positive relationship with economic, political, and cultural development [1]. As the world changes, governments have made increasing efforts to adjust their education policies in order to promote the individuals and communities [2, 3]. For example, the federal government of the United States published a Five-Year Strategic Plan for Science, Technology, Engineering and Mathematics (STEM) education in 2018, which emphasizes workforce development and technical skills [4]; moreover, the National Natural Science Foundation of China added a new application code F0701 in 2018 to support basic research in educational information science and technology [5]; Sweden also published a national strategy for the digitalization of the K-12 school system in 2017 with the aim of fostering adequate digital competence [6].

Recently, due to the threat of COVID-19, about 1.6 billion students (potentially exceeding 91% of total registered users, according to UNESCO) who would normally have been studying in school have been forced to be at home and 191 countries implemented country-wide school closures during the most difficult time [7]. In this context, online education, which plays an indispensable role in minimizing disruption to education, is developing on an unprecedented scale and gradually becoming a fashionable style of learning. Online education breaks free of the limitations of physical classrooms and enables to teach and learn flexibly anytime, anywhere [8, 9, 10]. At the same time, online learning has

the potential to bring huge educational benefits by providing each student with an optimal and adaptive learning experience. Online learning systems (such as Coursera, AS-SISTment and massive online open courses) [11, 12] have proven to be even more effective than traditional learning styles, owing to that they can offer more intelligent educational services, such as adaptive recommendations of individualized learning paths to students. In order to provide these intelligent services for each student, online learning systems continuously record a massive amount of available data about student-system interactions, which can be further mined to assess their knowledge levels and learning preferences. Concretely, Knowledge Tracing (KT) [13] is one of the most fundamental and critical research problems for intelligent education. KT utilizes a series of sequence modeling-oriented machine learning methods capable of exploiting educationally related data to monitor students' dynamic knowledge states. Nowadays, KT is widely applied in online intelligent tutoring systems and is also receiving growing attention[14, 15].

Fig. 1 presents a simple schematic diagram of the KT process. As the student works, the learning system continues to record the student's observational learning data, including exercises, the knowledge concepts (e.g., *equality, inequality, plane vector* and *probability*, which are represented in different colors in Fig. 1) contained in exercises, and students' answers (i.e., correct or incorrect responses). Benefiting from the development of intelligent education [16] and the methods of data anonymization [17], a large amount of side information is also recorded, such as response time, opportunity count and tutor intervention, which more completely reflect reflect students' learning behaviors. In addition, new technologies, enabling new approaches such as learning effective representations from limited educational data with crowdsourced labels, are also in constant development [18]. With reference to the recorded learning data, KT aims to maintain an estimate of student's changing knowledge state. Taking Fig. 1 as an intuitive example, the student's prior knowledge states on the four knowledge concepts are 0.2,

- Q. Liu, S. Shen, E. Chen (corresponding author) and Z. Huang are with the Anhui Province Key Laboratory of Big Data Analysis and Application, School of Data Science & School of Computer Science and Techonology, University of Science and Technology of China, Hefei, Anhui, 230026, China. Email: {qiliuql, huangzhy, cheneh}@ustc.edu.cn, {closer}@mail.ustc.edu.cn

- Y. Zheng is with the Research Institute of Science Education, Faculty of Education, Beijing Normal University, Haidian, Beijing, 100875, China. Email: {zhengyonghe}@bnu.edu.cn
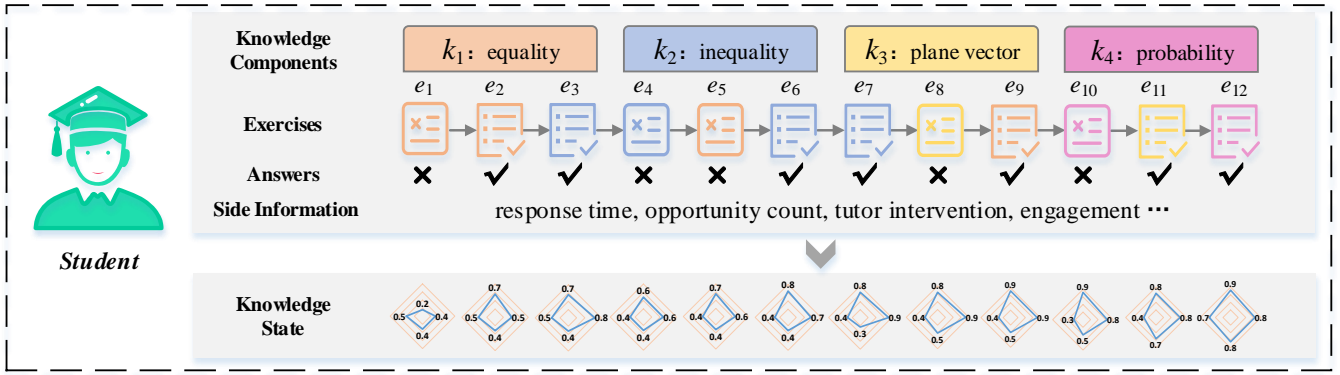
Fig. 1. A simple schematic diagram of knowledge tracing. Different knowledge concepts are represented in different colors, while exercises are also depicted in the color relevant to the knowledge concepts. During the learning process, different kinds of side information are also recorded in addition to the answers. The process by which the knowledge state changes is assessed by KT models and represented by the radar maps.

0.4, 0.4 and 0.5 respectively. While practicing, the student continues to absorb new knowledge, which can also be reflected by the gradually increased areas of the radar map that indicate the student's knowledge mastery. After a period of learning, the student's knowledge state reaches 0.9, 0.8, 0.8 and 0.7 on corresponding knowledge concepts, which indicates fairly good knowledge growth. KT models aim to monitor the process by which students' knowledge state changes in order to provide intelligent services to students, such as recommending exercises that suit their ability level. Indeed, KT is of great significance for both online learning systems and students. First, KT models enable the development of personalized adaptive learning systems. Once it has grasped the precise knowledge state of students, the learning system can customize more suitable learning schemes for different students, making it possible to teach students in accordance with their knowledge proficiency. Second, students themselves can also better understand their learning process, and they are encouraged to pay more attention to the knowledge concepts of which their mastery is poor in order to maximize their learning efficiency [19, 20].

Knowledge tracing has been studied for decades. Early KT-related studies can be traced back to the late 1970s; these works focused primarily on confirming the effectiveness of mastery learning [21]. Prior to 1990, the subject of KT was not discussed by researchers. To the best of our knowledge, Corbett and Anderson [13] were the first to introduce the concept of KT; these authors utilized Bayesian networks to model the student learning process, an approach referred to as Bayesian Knowledge Tracing. Subsequently, the significance of KT was recognized by a broader spectrum of people, and increasing attention has been channeled into KT-related research. For example, many logistic models have been applied for KT, including Learning Factor Analysis [22] and Performance Factor Analysis [23]. In recent years, deep learning has boosted research into the KT task, due to its powerful capacity to extract and represent features and its ability to discover intricate structure [24]. For example, Deep Knowledge Tracing introduced recurrent neural networks (RNNs) [25] into the KT task and was found to significantly outperform previous methods [26]. After that, by considering various characteristics of the learning sequence, many KT methods have introduced more types of neural networks to the KT task. For example, graph-based KT applied graph neural networks to model the knowledge

structure in knowledge concepts [27]. Some attention-based KT approaches have utilized the attention mechanism to capture dependencies between learning interactions [28, 29]. Moreover, due to the requirements of specific applications, many variants of KT models have also been continuously developed, and KT has already been successfully applied in many educational scenarios.

While novel KT models, as well as their massive variants and applications, continue to emerge, there have been few survey papers about this young research field, particularly as regards the emerging deep learning-based KT models. To this end, the present survey targets this research gap and aims to comprehensively review the current development and the state-of-the-art works on the KT task in a systematic manner. More specifically, as shown in Fig. 2, we first propose a new taxonomy of existing KT models from a technical perspective, which splits them into three categories: (1) probabilistic models, (2) logistic models and (3) deep learning-based models. Under this new taxonomy, we comprehensively review the basic KT models. We then introduce a large amount of variants of these basic KT models, which model more complete learning process in different learning phases. In addition, we present several typical applications of KT in different scenarios. Finally, we propose some potential future research directions. In general, this paper presents an extensive survey of knowledge tracing that can serve as basic guidelines for both researchers and practitioners in future research.

The remainder of this survey is organized as follows. Section 2 presents both the formal definition of the KT task and the taxonomy of KT models. Section 3 provides a review of the three categories of basic KT models. Section 4 describes the variants of basic KT models. Section 5 introduces the extensive applications of KT in different scenarios. Section 6 outlines some potential future research directions. Finally, section 7 summarizes the paper.

## 2 OVERVIEW

### 2.1 Problem Definition

In an online learning system, suppose there exist a set of students $\mathbb{S}$ and a set of exercises $\mathbb{E}$, where different students are asked to answer different exercises in order to achieve mastery of related knowledge. Each exercise is related to specific Knowledge Concepts (KCs; e.g.,
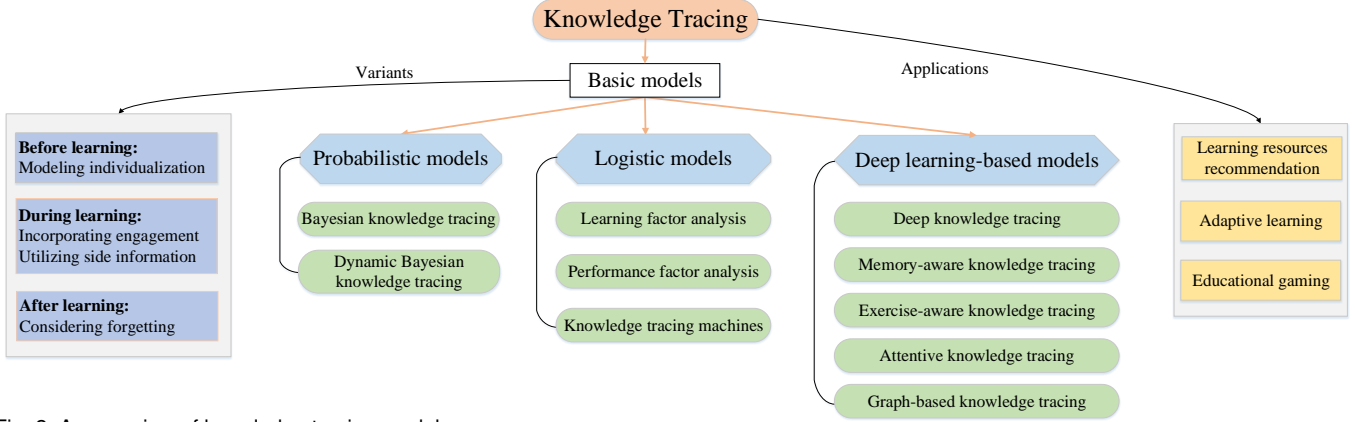
Fig. 2. An overview of knowledge tracing models.

*multiplication* and *friction*), which represents the basic unit of knowledge content. Generally speaking, the name given to the knowledge content differs across online learning platforms, for instance, it is named *skill* in ASSISTments [30]. To promote better understanding, we refer to these uniformly as knowledge concepts throughout this paper, and denote the set of all KCs as $\mathbb{KC}$. Besides, $M$ and $K$ are used to represent the total number of different exercises and KCs, respectively. The learning sequence of a student is represented as $\boldsymbol{X} = \{(e_1, a_1, r_1), (e_2, a_2, r_2), ..., (e_t, a_t, r_t), ..., (e_N, a_N, r_N)\}$, where the tuple $(e_t, a_t, r_t)$ represents the learning interaction of this student at the $t-$th time step, $e_t$ represents the exercise answered by the student, $a_t$ represents the correctness label (i.e., with 1 for correct and 0 for incorrect answers), $r_t$ stands for the side information recorded in this learning interaction, and $N$ is the sequence length of the learning interaction. The research problem of knowledge tracing can thus be defined as follows:

*Given the sequence of students' learning interactions in online learning systems, knowledge tracing aims to monitor students' changing knowledge states during the learning process and accurately predict their performance on future exercises; this information can be further applied to individualize students' learning schemes in order to maximize their learning efficiency.*

## 2.2 Categorization

As shown in Fig. 2, we divide and summarize the existing KT models according to their technical differences with a new taxonomy. In more detail, the proposed taxonomy splits existing KT methods into three categories: (1) probabilistic models, (2) logistic models, and (3) deep learning-based models. In addition to these basic KT models, we also introduce a large number of their variants, which respectively model more complete learning process from three various learning phases: (1) modeling individualization before learning, (2) incorporating engagement and utilizing side information during learning, and (3) considering forgetting after learning. Moreover, we also summarize the extensive applications of KT in different educational scenarios.

## 3 BASIC KNOWLEDGE TRACING MODELS

In this section, as shown in Table 1, we present the basic KT models according to our taxonomic framework. According

to the timeline of development, we will first introduce the probabilistic models, followed by the logistic models and finally the deep learning-based ones.

## 3.1 Probabilistic Models

The basic paradigm for probabilistic models in KT assumes that the learning process follows a Markov process, where students' latent knowledge state can be estimated by their observed learning performance [13]. In the following, we will present two basic probabilistic models in our taxonomy framework: the original Bayesian knowledge tracing (BKT) and the dynamic Bayesian knowledge tracing (DBKT).

### 3.1.1 Bayesian Knowledge Tracing

To the best of our knowledge, Bayesian Knowledge Tracing (BKT) was the first KT model to be proposed [13]. The topology of BKT's structure is illustrated in Fig. 3; here, the unshaded nodes represent unobservable latent knowledge states, while the shaded nodes represent the observable answers of the student.

In fact, BKT is a special case of Hidden Markov Model (HMM). There are two types of parameters in HMM: transition probabilities and emission probabilities. In BKT, the transition probabilities are defined by two learning parameters: (1) $P(T)$, the probability of transition from the unlearned state to the learned state; (2) $P(F)$, the probability of forgetting a previously known KC, which is assumed to be zero in BKT. Moreover, the emission probabilities are determined by two performance parameters: (1) $P(G)$, the probability that a student will guess correctly in spite of non-mastery; (2) $P(S)$, the probability a student will make a mistake in spite of mastery. Furthermore, the parameter $P(L_0)$ represents the initial probability of mastery. BKT assumes a two-state student modeling framework: a knowledge concept is either learned or unlearned by the student, and there is no forgetting once a student has learned the knowledge. Given the observations of the student's learning interactions, the following equation is used to estimate the knowledge state and the probability of correct answers:

$$P(L_n) = P(L_n|Answer) + (1 - P(L_n|Answer))P(T),$$
$$P(C_{n+1}) = P(L_n)(1 - P(S)) + (1 - P(L_n))P(G), \quad (1)$$

where $P(L_n)$ is the probability that a KC is mastered by the student at the $n$-th learning interaction, $P(C_{n+1})$ is the

TABLE 1
A summary of different types of basic knowledge tracing models.

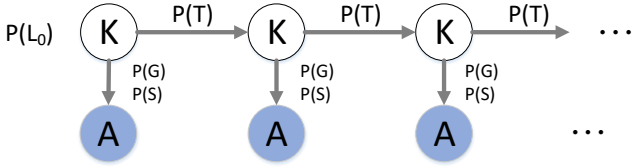| Category | Typical approach | Technique | KC relationship | Knowledge state |
|---|---|---|---|---|
| Probabilistic models | Bayesian knowledge tracing [13] | Bayesian networks | independent | unobservable node in HMM |
| | dynamic Bayesian knowledge tracing [31] | Dynamic Bayesian networks | pre-defined | |
| Logistic models | learning factor analysis [22] | logistic regression | independent | the output of logistic regression function |
| | performance factor analysis [23] | | | |
| | knowledge tracing machines [32] | factorization machines | | |
| Deep learning-based models | deep knowledge tracing [26] | RNN/LSTM | discover automatically | the hidden state |
| | memory-aware knowledge tracing [33] | memory networks | correlation weights | *value* matrix |
| | exercise-aware knowledge tracing [34, 35] | semantic exercise representations | related coefficients | weighted sum of historical states |
| | attentive knowledge tracing [28, 29] | self-attention mechanism | attention weights | attentive historical knowledge state |
| | graph-based knowledge tracing [27] | graph neural networks | edges in graph | aggregate in the graph |



Fig. 3. The topology of Bayesian Knowledge Tracing [13]. $K$ are the unobserved knowledge nodes, $A$ are the observed performance (answer) nodes, $P(L_0)$ represents the initial probability, $P(T)$ is the transition probability, $P(G)$ is the guessing probability and $P(S)$ is the slipping probability.
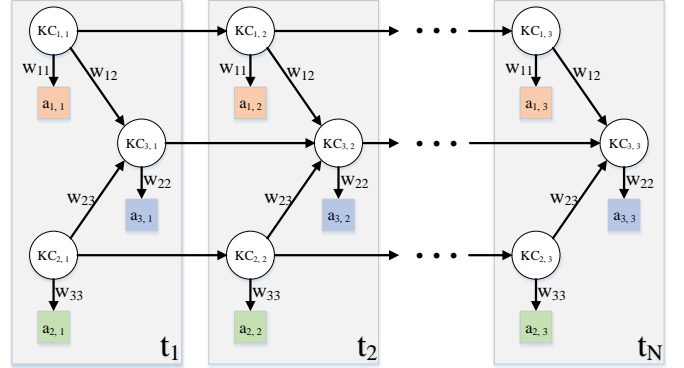


Fig. 4. The topology of Dynamic Bayesian Knowledge Tracing [31]. $KC_{i,j}$ denotes the knowledge state state of KC $i$ at time step $j$, while $a_{i,j}$ are the corresponding observed answer nodes.

probability of correct answers at the next learning interaction. $P(L_n)$ is the sum of two probabilities: (1) the probability that the KC is already in a learned state, which is not affected by the answer; (2) the probability that the student's knowledge state of the KC will make the transition to the learned state. The posterior probability $P(L_n|Answer)$ is estimated by a Bayesian inference scheme, as follows:

$$P(L_n|correct) = \frac{P(L_{n-1})(1 - P(S))}{P(L_{n-1})(1 - P(S)) + (1 - P(L_{n-1}))P(G)},$$
$$P(L_n|incorrect) = \frac{P(L_{n-1})P(S)}{P(L_{n-1})P(S) + (1 - P(L_{n-1}))(1 - P(G))}. \quad (2)$$

### 3.1.2 Dynamic Bayesian Knowledge Tracing

BKT models the parameters of each KC individually, i.e., it employs one BKT model for each KC. However, KCs are not completely independent of each other, but rather hierarchical and closely related [36]. Dynamic Bayesian networks are able to jointly represent multiple skills within one model, which can potentially increase the representational power of knowledge tracing. Therefore, Käser et al. [31] proposed dynamic Bayesian knowledge tracing (DBKT) to model the prerequisite hierarchies and relationships within KCs based on dynamic Bayesian networks, which considers different KCs jointly within a single model.

In DBKT, a student's knowledge mastery is also represented by binary latent variables and can be estimated from the student's learning interactions. Fig. 4 depicts the topology of DBKT; here, the unshaded circular node $KC_i$ represents the binary mastery variable of the KC, while $i$ is the index of KC. The shaded rectangular nodes represent observable variables (i.e., binary student answers) associated with the corresponding KCs. In contrast to BKT, DBKT models the dependencies between the different KCs; for example, $KC_1$ and $KC_2$ are prerequisites for mastering

$KC_3$ in Fig. 4, therefore students' mastery of $KC_3$ depends on their mastery of $KC_1$ and $KC_2$.

Let $H$ denote the unobserved variables, i.e., lack of student answers and binary mastery variables. In Fig. 4, suppose that the student is answering an exercise associated with $KC_1$ at time step $t_1$ correctly, i.e., $a_{1,1} = 1$. The set of observed variables is then $a_m = a_{1,1}$ and the set of unobserved variables is $h_m = \{KC_{1,1}, KC_{2,1}, KC_{3,1}, a_{2,1}, a_{3,1}\}$. The objective of DBKT is to find the parameters $\theta$ that maximize the likelihood of the joint probability $p(a_m, h_m|\theta)$. The log-likelihood can alternatively be formulated using a log-linear model, as follows:

$$L(\boldsymbol{w}) = \sum_m ln(\sum_{h_m} exp(\boldsymbol{w}^T \Phi(a_m, h_m) - ln(Z))), \quad (3)$$

where $\Phi : A \times H \to \mathbb{R}^F$ denote a mapping from the observed space $A$ and the latent space $H$ to an $F$-dimensional feature vector. Moreover, $Z$ is a normalizing constant, and $\boldsymbol{w}$ denotes the weights of the model.

## 3.2 Logistic Models

In contrast to probabilistic models, logistic models are a large class of models based on logistic functions, the underlying concept behind which is that the probability of answering exercises correctly can be represented by a mathematical function of student and KC parameters. In logistic models, the binary answers (*correct/incorrect*) submitted the students follow a Bernoulli distribution. These models first take advantage of different factors in student's learning

interactions to compute an estimation of the student and KC parameters, then utilize a logistic function to transform this estimation into the prediction of the probability of mastery [37]. We will introduce three logistic models in this section: (1) learning factor analysis (LFA), (2) performance factor analysis (PFA) and (3) knowledge tracing machines (KTM).

### 3.2.1 Learning Factor Analysis

The LFA model [22] considers the following learning factors:

- Initial knowledge state: parameter $\alpha$ estimates the initial knowledge state of each student;
- Easiness of KCs: some KCs are more likely to be known, and parameter $\beta$ captures the easiness of different KCs;
- Learning rate of KCs: some KCs are easier to learn than others, and parameter $\gamma$ denotes the learning rate of KCs.

The learning rates of different students are assumed to be the same. The standard form of the LFA model is as follows:

$$
\theta = \sum_{i \in N} \alpha_i S_i + \sum_{j \in KCs} (\beta_j + \gamma_j T_j) K_j,
$$
$$
p(\theta) = \frac{1}{1 + e^{-\theta}},
$$
(4)

where $S_i$ is the covariates for the student $i$, $T_j$ represents the covariate for the number of practice opportunities on KC $j$, $K_j$ is the covariate for KC $j$, $\theta$ is the estimation of the probability of student and KC parameters, and $p(\theta)$ is the estimation of the probability of a correct answer.

### 3.2.2 Performance Factor Analysis

The PFA model [23] can be seen as an extension of the LFA model that is especially sensitive to the strongest indicator of student learning performance. In contrast to the LFA model, the PFA model considers different factors, namely the following:

- Previous failures: parameter $f$ is the prior failures for the KC of the student;
- Previous successes: parameter $s$ represents the prior successes for the KC of the student;
- Easiness of KCs: parameter $\beta$ menas the easiness of different KCs, which is the same as in the LFA model.

The standard PFA model takes the following form:

$$
\theta = \sum_{j \in KCs} (\beta_j + \mu_j s_{ij} + \nu_j f_{ij}),
$$
$$
p(\theta) = \frac{1}{1 + e^{-\theta}},
$$
(5)

where $\mu$ and $\nu$ are the coefficients for $s$ and $f$, which denote the learning rates for successes and failures, respectively.

### 3.2.3 Knowledge Tracing Machines

The KTM model [32] takes advantage of factorization machines (FMs) [38, 39] to generalize previous logistic models to higher dimensions. FMs were originally proposed as a general predictor that works with any real valued feature vector, which can model all interactions between variables using factorized parameters [40]. FMs provide a means of encoding side information about exercises or students into the model; thus, KTM is able to make use of all the information available at hand.
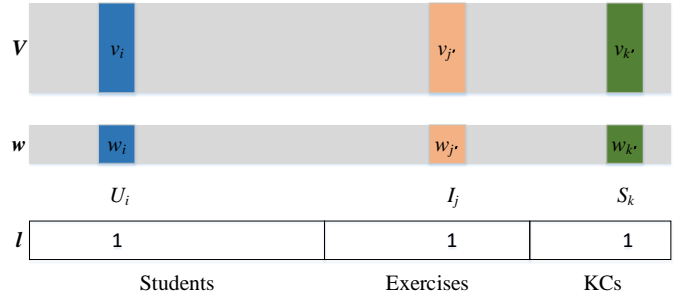


Fig. 5. Example of activation of a knowledge tracing machine [32]. $V$ refers to the matrix of embeddings, $w$ refers to the vector of biases, and $l$ is the encoding vector of the learning interaction.

Fig. 5 presents an illustrative example of KTM, which models the knowledge mastery of the student based on a sparse set of weights for all features involved in the event. Let $L$ be the number of features; here, the features can be related either to students, exercises, KCs or any other side information. The learning interaction is encoded by a sparse vector $\boldsymbol{l}$ of length $L$, where $l_i > 0$ if feature $i$ is involved in the interaction. The probability $p(\theta)$ of the correct answer is determined by the following equations:

$$
\theta = \mu + \sum_{i=1}^{L} w_i l_i + \sum_{1 \le i < j \le L} l_i l_j \langle \boldsymbol{v_i}, \boldsymbol{v_j} \rangle,
$$
$$
p(\theta) = \frac{1}{1 + e^{-\theta}},
$$
(6)

where $\mu$ is the global bias, and the feature $i$ is modeled by a bias $w_i \in \boldsymbol{R}$ and an embedding $\boldsymbol{v_i} \in \boldsymbol{R}^d$; here, $d$ is the dimension. Note that only features with $l_i > 0$ will have impacts on the predictions.

## 3.3 Deep Learning-based Models

Notably, the cognitive process can be influenced by many factors at both the macro and micro level. It is difficult for probabilistic or logistic models to adequately capture a cognitive process of this complexity [26]. Deep learning has a powerful ability to achieve non-linearity and feature extraction, making it well suited to modeling the complex learning process, especially when a much larger amount of learning interaction data is available [41]. In recent years, many research works on deep learning-based KT models have been proposed and achieved quite good performance. Nevertheless, deep learning-based models are poorly interpretable due to their end-to-end learning strategy, which limits their further applicability owing to the crucial significance of interpretability for students modeling. According to our taxonomy framework, we will introduce deep learning-based models from five aspects based on the technical differences: (1) deep knowledge tracing, (2) memory-aware knowledge tracing, (3) exercise-aware knowledge tracing, (4) attentive knowledge tracing, and (5) graph-based knowledge tracing.

### 3.3.1 Deep Knowledge Tracing

Deep knowledge tracing (DKT), which is the first approach to introduce deep learning into knowledge tracing [26], utilizes recurrent neural networks (RNNs) [25] to model the students' knowledge states. DKT applies RNNs to process
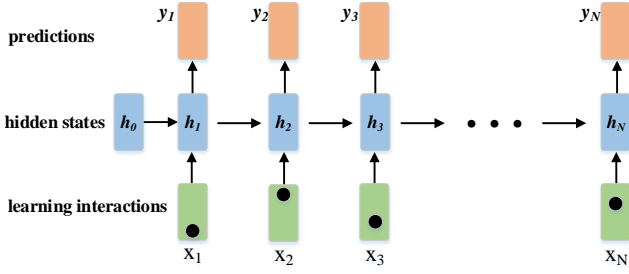
Fig. 6. Deep Knowledge Tracing [26]. $x_i$ are the input embeddings of students' learning interactions, $h_i$ are the hidden states that represent students' knowledge state, and $y_i$ are the predicted answers.

the input sequence of learning interactions over time, maintaining a hidden state that implicitly contains information about the history of all past elements of the sequence. The hidden state evolves based on both the previous knowledge state and the present input learning interaction [26]. DKT provides a high-dimensional and continuous representation of the knowledge state, making them better able to model the complex learning process. Generally, RNNs' variant long short-term memory (LSTM) networks [42] are more commonly used in the implementation of DKT, which is made more powerful through considering forgetting.

Fig. 6 illustrates the process of deep knowledge tracing. In DKT, exercises are represented by their contained KCs. For datasets with different number of KCs, DKT applies two different methods to convert students' learning interactions $\boldsymbol{X} = \{(e_1,a_1),(e_2,a_2),...,(e_t,a_t),...,(e_N,a_N)\}$ into a sequence of fixed-length input vectors. More specifically, for datasets with a small number $K$ of unique KCs, $\boldsymbol{x}_t \in \{0,1\}^{2K}$ is set as a one-hot embedding, where $\boldsymbol{x}_t^k = 1$ if the answer $a_t$ of the exercise with KC $k$ was correct or $\boldsymbol{x}_t^{k+K} = 1$ if the answer was incorrect. For datasets with a large number of unique KCs, one-hot embeddings are too sparse; therefore, DKT sets each input vector $\boldsymbol{x}_t$ to a corresponding random vector, then takes the embedded learning sequence as the input of RNNs and applies a linear mapping and activation function to the output hidden states to obtain the knowledge state of students:

$$\begin{aligned}
\boldsymbol{h}_t &= tanh(\boldsymbol{W}_{hs}\boldsymbol{x}_t + \boldsymbol{W}_{hh}\boldsymbol{h}_{t-1} + \boldsymbol{b}_h), \\
\boldsymbol{y}_t &= \sigma(\boldsymbol{W}_{yh}\boldsymbol{h}_t + \boldsymbol{b}_y),
\end{aligned} \tag{7}$$

where $\sigma$ is the sigmoid function, $tanh$ is the activation function, $\boldsymbol{W}_{hs}$ is the input weight parameter, $\boldsymbol{W}_{hh}$ is the recurrent weight parameter, $\boldsymbol{W}_{yh}$ is the readout weight parameter, $\boldsymbol{b}_h$ and $\boldsymbol{b}_y$ are the bias terms, respectively.

Benefiting from the advantages of RNNs, DKT has demonstrated a superior performance contrast to the probabilistic and logistic models. Moreover, DKT can further discover the exercise relationships after training. Nevertheless, DKT also has certain shortcomings that cannot be ignored. For example, the DKT model lacks interpretability: it is difficult to figure out how the hidden states can represent students' knowledge state, and it cannot explicitly determine a student's level of knowledge mastery from the hidden state [41]. Yeung and Yeung [43] further revealed that there are two unreasonable phenomena in DKT that violate the common sense, i.e., (1) it fails to reconstruct the observed input, and (2) the predicted knowledge state is not consistent across time-steps. Overall, DKT remains a promising KT model [44].
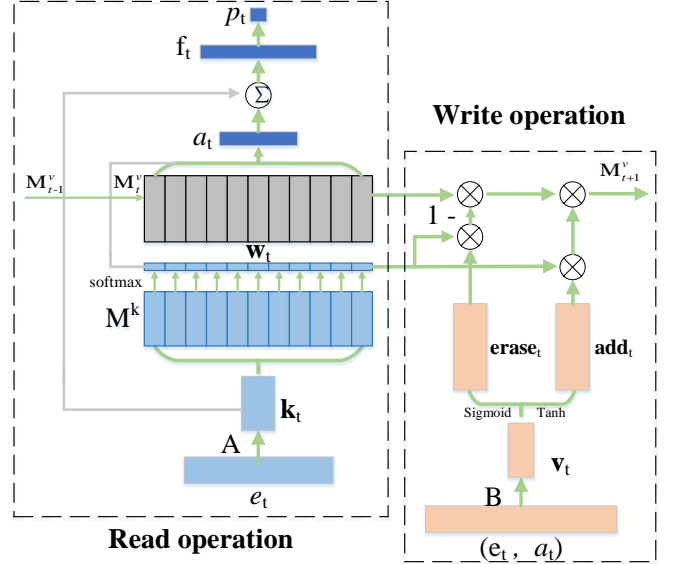


Fig. 7. The architecture for dynamic key-value memory networks [33].

### 3.3.2 Memory-aware Knowledge Tracing

In order to enhance the interpretability of deep learning-based KT models, memory-aware knowledge tracing introduces an external memory module [45] to store the knowledge concepts and update the corresponding knowledge mastery of the student. The most representative of these models is Dynamic Key-Value Memory Networks (DKVMN) for knowledge tracing [33], which can accurately point out students' specific knowledge state on KCs. DKVMN initializes a static matrix called a $key$ matrix to store latent KCs and a dynamic matrix called a $value$ matrix to store and update the mastery of corresponding KCs through read and write operations over time.
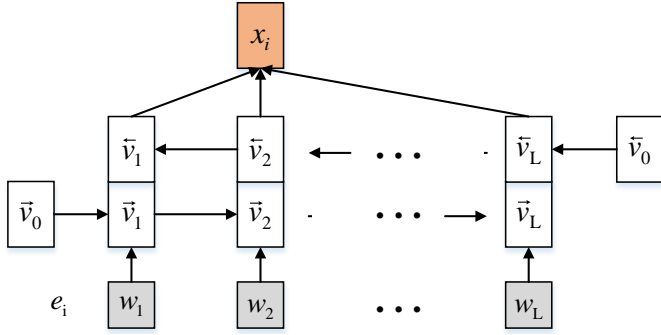
As shown in Fig. 7, an embedding matrix is first defined to obtain the embedding vector $k_t$ of the exercises. A correlation weight $\boldsymbol{w}_t$ is then obtained by taking the inner product between the exercise embedding $k_t$ and the $key$ vectors $M^k$, following by the softmax activation:

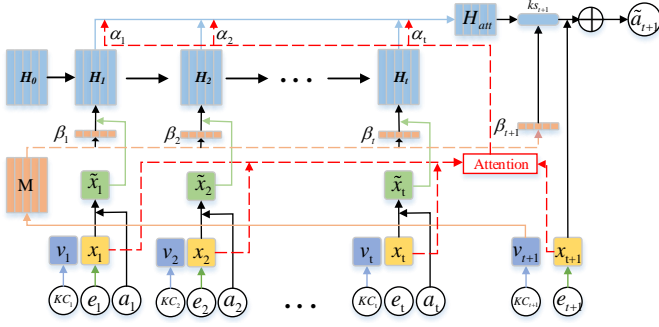$$\boldsymbol{w}_t = Softmax(k_t M^k), \tag{8}$$

where the correlation weight $\boldsymbol{w}_t$ represents the correlation between the exercises and all latent KCs. Both the read and write operations will use $\boldsymbol{w}_t$.

For the read operation, DKVMN can predict student performance based on the student's knowledge mastery. More specifically, DKVMN reads students' mastery of the exercise $\boldsymbol{r}_t$ with reference to the weighted sum of all memory vectors in the $value$ matrix using the correlation weight. Then the read content and the input exercise embeddings are then concatenated together and passed to a fully connected layer to yield a summary vector $\boldsymbol{f}_t$, which contains both the student's knowledge mastery and the prior difficulty of the exercise. Furthermore, the student's performance can be predicted by applying another fully connected layer with a sigmoid activation function to the summary vector:

$$\begin{aligned}
\boldsymbol{r}_t &= \sum_{i=1}^{N} w_t(i) M_t^v(i), \\
\boldsymbol{f}_t &= tanh(\boldsymbol{W}_f[\boldsymbol{r}_t, k_t] + \boldsymbol{b}_f), \\
p_t &= \sigma(\boldsymbol{W}_p \boldsymbol{f}_t + \boldsymbol{b}_p),
\end{aligned} \tag{9}$$

(a) The process of learning the semantic representation of exercises.



(b) EKT with attention mechanism.

Fig. 8. Exercise-aware Knowledge Tracing [34, 35].

where $\boldsymbol{W}_f$ and $\boldsymbol{W}_p$ are the weight parameters, while $\boldsymbol{b}_f$ and $\boldsymbol{b}_p$ are bias terms.

For the write operation, after an exercise has been responded to, DKVMN updates students' knowledge mastery (i.e., the $value$ matrix) based on their performance. In more detail, the learning interaction $(e_t, a_t)$ is first embedded with an embedding matrix $B$ to obtain the student's knowledge growth $\boldsymbol{v}_t$. Then DKVMN calculates an erase vector $\boldsymbol{erase}_t$ from $\boldsymbol{v}_t$ and decides to erase the previous memory $value$ with reference to both the erase vector and the correlation weight $\boldsymbol{w}_t$. Following erasure, the new memory vectors are updated by the new knowledge state and the add vector $\boldsymbol{add}_t$, which forms an $erase$-followed-by-$add$ mechanism that allows forgetting and strengthening knowledge mastery in the student's learning process:

$$
\begin{aligned}
\boldsymbol{erase}_t &= \sigma(\boldsymbol{W}_e \boldsymbol{v}_t + \boldsymbol{b}_e), \\
\widetilde{M}_t^v(i) &= M_{t-1}^v(i)[1 - w_t(i)\boldsymbol{erase}_t], \\
\boldsymbol{add}_t &= tanh(\boldsymbol{W}_d \boldsymbol{v}_t + \boldsymbol{b}_d), \\
M_t^v(i) &= \widetilde{M}_t^v(i) + w_t(i)\boldsymbol{add}_t,
\end{aligned}
\tag{10}
$$

where $\boldsymbol{W}_e$ and $\boldsymbol{W}_d$ are the weight parameters, while $\boldsymbol{b}_e$ and $\boldsymbol{b}_d$ are bias terms.

Abdelrahman and Wang [46] pointed out that although DKVMN modeled students' knowledge state through their most recent practices, it failed to capture long-term dependencies in learning process. Therefore, they proposed a Sequential Key-Value Memory Network (SKVMN) to combine the strengths of DKT's recurrent modelling capacity and DKVMN's memory capacity for knowledge tracing. In SKVMN, a modified LSTM called *Hop-LSTM* is used to hop across LSTM cells according to the relevance of the latent KCs, which directly captures the long-term dependencies. In the write process, when calculating the knowledge growth

of a new exercise, SKVMN enables it to consider the current knowledge state in order to get more reasonable results.

### 3.3.3 Exercise-aware Knowledge Tracing

In online education systems, the text content is of great significance for the understanding of exercises (e.g., similarity and difficulty). For example, Huang et al. [47] used text materials to automatically predict the difficulty of exercises and Liu et al. [48] utilized the text content to find similar exercises. Yin et al. [49] further proposed a pre-training model called QuesNet for comprehensively learning the unified representations of heterogeneous exercises. In fact, exercises' text contents also have significant impacts on students' performance, as understanding exercises is the first step to answer them. Therefore, Liu et al. [35] proposed the Exercise-aware Knowledge Tracing (EKT) model to leverage the effectiveness of the text content of exercises in order to enhance the KT process. More specifically, instead of using one-hot encoding of exercises, EKT automatically learns the semantic representation $x_i$ of each exercise from its text content $e_i$. As shown in Fig. 8(a), EKT first uses $Word2vec$ [50] to train a pre-trained word embedding vector for each word $w_i$ in exercise $e_i$, then it constructs a bidirectional LSTM, which capturing the word sequence from both forward and backward directions, to learn the semantic word representation. Finally, the element-wise max-pooling operation is utilized to merge $L$ words' contextual representations into a global embedding $x_i$ as $x_i = max(v_1, v_2, ..., v_L)$. After obtaining the semantic representation $x_i$ of each exercise, in order to distinguish the different influences of correct and incorrect answers on a student's knowledge state, EKT extends the answer $a_t$ to a feature vector $\mathbf{0} = (0, 0, ..., 0)$ with the same dimensions as $x_t$ and represents the learning interaction $\widetilde{x}_t$ as follows:

$$
\widetilde{x}_t = \begin{cases} [x_t \oplus \mathbf{0}], & \text{if} \quad a_t = 1, \\ [\mathbf{0} \oplus x_t], & \text{if} \quad a_t = 0. \end{cases}
\tag{11}
$$

In order to explicitly measure the extent to which the student has mastered a certain KC, EKT further incorporates the information of KCs associated with each exercise. As shown in Fig. 8(b), a memory module consisted of a matrix $M$ is set up to represent KCs. The KCs of each exercise is converted into a one-hot encoding $KC_t \in \{0, 1\}^K$ with dimension equaling to the total number K of all KCs. An embedding matrix $\boldsymbol{W}_K \in \mathbb{R}^{K \times d_k}$ transfers $KC_t$ into a low-dimensional vector $v_t \in \mathbb{R}^{d_k}$ as follows: $v_t = \boldsymbol{W}_K^T KC_t$. Another static memory network [45] is then utilized to calculate the knowledge impact $\beta_t^i$, which quantifies the correlation weights between the $i$-th KC of the exercise and each knowledge memory vector in $M$, as follows:

$$
\beta_t^i = Softmax(v_t^T) = \frac{\exp(v_t^T \boldsymbol{M}_i)}{\sum_{i=1}^K \exp(v_t^T \boldsymbol{M}_i)}.
\tag{12}
$$

With the knowledge impact $\beta_t$ of each exercise, the input $\widetilde{x}_t$ is replaced by a new joint representation: $\widetilde{x}_t^i = \beta_t^i \widetilde{x}_t$. At

each learning step $t$, EKT updates the student's knowledge state $H_t^i \in \mathbb{R}^{K \times d_h}$ by the LSTM networks [42], as follows:

$$
\begin{aligned}
i_t &= \sigma(\boldsymbol{W}_{\widetilde{x}i}\widetilde{x}_t + \boldsymbol{W}_{Hi}H_{t-1}^i + \boldsymbol{b}_i), \\
f_t &= \sigma(\boldsymbol{W}_{\widetilde{x}f}\widetilde{x}_t + \boldsymbol{W}_{Hf}H_{t-1}^i + \boldsymbol{b}_f), \\
o_t &= \sigma(\boldsymbol{W}_{\widetilde{x}o}\widetilde{x}_t + \boldsymbol{W}_{Hc}H_{t-1}^i + \boldsymbol{b}_o), \\
c_t &= f_t \cdot c_{t-1} + i_t \cdot tanh(\boldsymbol{W}_{\widetilde{x}c}\widetilde{x}_t + \boldsymbol{W}_{Hc}H_{t-1}^i + \boldsymbol{b}_c)), \\
H_t^i &= o_t \cdot tanh(c_t),
\end{aligned}
\tag{13}
$$

where $i_t$ is the input gate, $f_t$ is the forget gate and $o_t$ is the output gate; moreover, $\boldsymbol{W}_{\widetilde{x}*}$ and $\boldsymbol{W}_{H*}$ are the weight parameters, while $\boldsymbol{b}_*$ are the bias terms, respectively.

Next, EKT utilizes the attention mechanism to enhance the effect of student performance on similar exercises in order to predict future student performance. The matrix $H_{att}$ represents the attentive hidden state, where each slot $H_{att}^i$ can be computed as follows:

$$
H_{att}^i = \sum_{j=1}^{T} \alpha_j H_j^i, \quad \alpha_j = cos(x_{t+1}, x_t).
\tag{14}
$$

Finally, EKT can predict the student's performance on the next exercise $e_{t+1}$ as follows:

$$
\begin{aligned}
s_{t+1} &= \sum_{i=1}^{K} \beta_{t+1}^i H_{att}^i, \\
y_{t+1} &= ReLU(\boldsymbol{W}_1 \cdot [s_{t+1} \oplus x_{t+1}] + \boldsymbol{b}_1), \\
a_{t+1} &= \sigma(\boldsymbol{W}_2 \cdot y_{t+1} + \boldsymbol{b}_2),
\end{aligned}
\tag{15}
$$

where $\boldsymbol{W}_1$ and $\boldsymbol{W}_2$ are weight parameters, $\boldsymbol{b}_1$ and $\boldsymbol{b}_2$ are bias terms, and $\oplus$ is the concatenation operation.

In addition to exploring the text content in order to capture meaningful exercise information, Liu et al. [51] presented a method for obtaining pre-trained exercise embeddings so as to improve the accuracy of knowledge tracing. In fact, in KT scenarios, the explicit exercise-KC relations and the implicit exercise similarity as well as KC similarity exists simultaneously. To capture all these relations in exercise embeddings, Liu et al. [51] first represented them together with exercise difficulties as a bipartite graph. They then utilized a product layer to fuse these features in the defined bipartite graph to obtain the pre-trained exercise embeddings. Finally, their extensive experiments indicated that the obtained exercise embeddings can significantly improve the performance of some KT models, such as DKT.

### 3.3.4 Attentive Knowledge Tracing

In the development of deep learning, transformer was the first approach to be proposed for neural machine translation [52]; this model abandons recurrence and relies entirely on an attention mechanism to capture global dependencies within a sequence. The transformer has demonstrated superior power in feature extraction and dependency capture while maintaining high computational efficiency. Some representative transformer-based pre-training models, such as BERT [53], have obtained state-of-the-art results on many natural language processing tasks. Pandey and Karypis [28] proposed a self-attentive model for knowledge tracing (SAKT), which directly applied the transformer to capture long-term dependencies between students' learning interaction with few modifications, and achieved fairly good performance. Moreover, Wang et al. [54] proposed an adaptive sparse self-attention network to generate the missing features and simultaneously obtained fine-grained prediction of student performance. Zhu et al. [55] found there was a vibration problem in DKT and presented an attention-based knowledge tracing model to solve it, which also further used Finite State Automaton (FSA) to provide deep analysis about knowledge state transition.

However, the complexity of the KT task, which is caused by the interactions between students and exercises, limits the performance of the simple transformer application. Therefore, Choi et al. [56] proposed a model named separated self-attentive neural knowledge tracing (SAINT) to improve self-attentive computation for knowledge tracing adaptation. More specifically, SAINT has an encoder-decoder structure, where the exercise and answer embeddings are separately encoded and decoded by self-attention layers. This separation of input allows SAINT to stack self-attention layers multiple times and capture complex relations among exercises and answers. Subsequently, the SAINT+ model [57] is proposed to incorporate two temporal features into SAINT: namely, the answering time for each exercise and the interval time between two continuous learning interactions. Both SAINT and SAINT+ have achieved superior performance relative to SAKT on the EdNet dataset [58], one of the largest publicly available datasets for educational data mining.

Similarly, Ghosh et al. [29] observed that SAKT did not outperform DKT and DKVMN in their experiments. Unlike SAINT and SAINT+, these authors presented a context-aware attentive knowledge tracing (AKT) model, incorporating self-attention mechanism with cognitive and psychometric models. AKT comprises four modules: Rasch model-based embeddings, exercise encoder, knowledge encoder and knowledge retriever. More specifically, the classic and powerful Rasch model in psychometrics [59] was utilized to construct embeddings for exercises and KCs. The embedding of the exercise $e_t$ with KC $c_t$ is constructed as follows:

$$
\boldsymbol{x}_t = \boldsymbol{c}_{c_t} + \mu_{e_t} \cdot \boldsymbol{d}_{c_t},
\tag{16}
$$

where $\boldsymbol{c}_{c_t} \in \mathbb{R}^D$ is the embedding of the KC of this exercise, $\boldsymbol{d}_{c_t} \in \mathbb{R}^D$ is a vector that summarizes the variation in exercises with the related KC, and $\mu_{e_t} \in \mathbb{R}^D$ is a scalar difficulty parameter that controls the extent to which this exercise deviates from the related KC. The exercise-answer tuple $(e_t, a_t)$ is similarly extended using the scalar difficulty parameter for each pair:

$$
\boldsymbol{y}_t = \boldsymbol{q}_{(c_t, a_t)} + \mu_{e_t} \cdot \boldsymbol{f}_{(c_t, a_t)},
\tag{17}
$$

where $\boldsymbol{q}_{(c_t, a_t)} \in \mathbb{R}^D$ and $\boldsymbol{f}_{(c_t, a_t)} \in \mathbb{R}^D$ are KC-answer embedding and variation vectors. Through such embedding, exercises labeled as the same KCs are determined to be closely related while keeping important individual characteristics. Then, the input of the exercise encoder is the exercise embeddings $\{\boldsymbol{e}_1, ..., \boldsymbol{e}_t\}$ and the output is a sequence of context-aware exercise embeddings $\{\widetilde{\boldsymbol{e}}_1, ..., \widetilde{\boldsymbol{e}}_t\}$. AKT designs a monotonic attention mechanism to accomplish the above process, where the context-aware embedding of each

exercise depends on both itself and the previous exercises, i.e., $\widetilde{e}_t = f_{enc_1}(e_1, ..., e_t)$. Similarly, the knowledge encoder takes exercise-answer embeddings $\{y_1, ..., y_t\}$ as input and outputs a sequence of context-aware embeddings of the knowledge acquisitions $\{\widetilde{y}_1, ..., \widetilde{y}_t\}$ using the same monotonic attention mechanism, which are also determined by students' answers to both the current exercise and prior exercises, i.e., $\widetilde{y}_t = f_{enc_1}(y_1, ..., y_t)$. Finally, the knowledge retriever takes the context-aware exercise embedding $\widetilde{e}_{1:t}$ and exercise-answer pair embeddings $\widetilde{y}_{1:t}$ as input and outputs a retrieved knowledge state $h_t$ for the current exercise. Since the student's current knowledge state depends on the related answering exercise, it is also context-aware in AKT. The novel monotonic attention mechanism proposed in AKT is based on the assumptions that the learning process is temporal and students' knowledge will decay over time. Therefore, the scaled inner-product attention mechanism utilized in the original transformer is not suitable for the KT task. AKT uses exponential decay and a context-aware relative distance measure to computes the attention weights. Finally, AKT achieves outstanding performance on predicting students' future answers, as well as demonstrating interpretability due to the combination of cognitive and psychometric models.

Besides, as the contextual information of exercises has been identified as highly significant in EKT, Pandey and Srivastava [60] proposed a relation-aware self-attention model for knowledge tracing (RKT), which utilizes contextual information to enhance the self-attention mechanism. RKT defines a concept called the relation coefficients to capture the relations between exercises, which are obtained from modeling the textual content of the exercises and students' forgetting behaviors respectively. The contextual exercise representation is then fed to the self-attention layer to trace students' knowledge state.

### 3.3.5 Graph-based Knowledge Tracing

Graph neural networks (GNNs), which are designed to handle the complex graph-related data, have developed rapidly in recent years [61]. In this case, the graph represents a kind of data structure that models a set of objects (nodes) and their relationships (edges). From a data structure perspective, there is a naturally existing graph structure within the KCs. Therefore, incorporating the graph structure of the knowledge concepts as a relational inductive bias should be able to improve the performance of knowledge tracing models. Nakagawa et al. [27] presented the graph-based knowledge tracing (GKT), which conceptualizes the potential graph structure of the knowledge concepts as a graph $G = (V, E)$; here, nodes $V = \{v_1, v_2, ..., v_N\}$ represents the set of KCs and the edges $E \subseteq V \times V$ represents the relationships between these KCs, $h^t = \{h_{i \in V}^t\}$ represents the student's temporal knowledge state after answering the exercise at time $t$. The architecture for graph-based knowledge tracing is presented in Fig. 9, which is composed of three parts: (1) *aggregate*, (2) *update* and (3) *predict*.

In the *aggregate* module, GKT aggregates the temporal knowledge state and the embedding for the answered KC $i$ and its neighboring KC $j$:

$$h_k'^t = \begin{cases} [h_k^t, a^t E_s] & (k = i), \\ [h_k^t, E_e(k)] & (k \neq i), \end{cases} \quad (18)$$


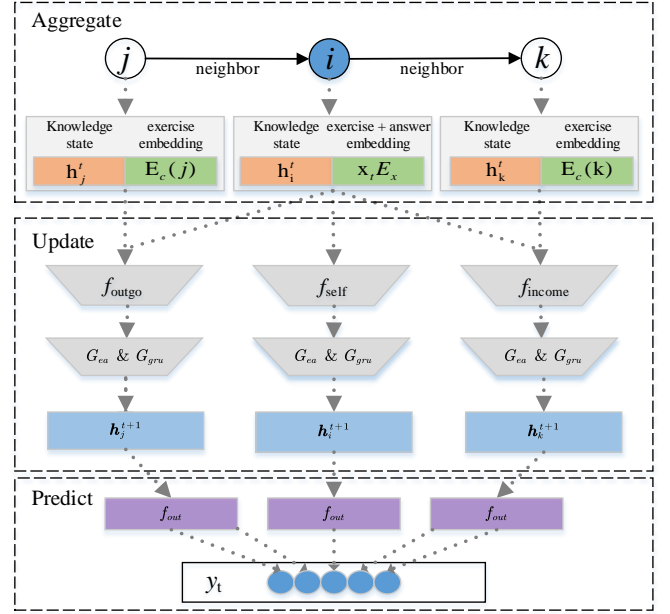
Fig. 9. The architecture for graph-based knowledge tracing [27].

where $a^t$ represents the exercises answered correctly or incorrectly at time step $t$, $E_s$ is the embedding matrix for the learning interactions, $E_e$ is the embedding matrix for the KC, and $k$ represents the $k$-th row of $E_e$.

In the *update* module, GKT updates the temporal knowledge state based on the aggregated features and the knowledge graph structure, as follows:

$$\begin{aligned} m_k^{t+1} &= \begin{cases} f_{self}(h_k'^t)(k = i), \\ f_{neighbor}(h_i'^t, h_k'^t)(k \neq i), \end{cases} \\ \widetilde{m}_k^{t+1} &= G_{ea}(m_k^{t+1}), \\ h_k^{t+1} &= G_{gru}(\widetilde{m}_k^{t+1}, h_k^t), \end{aligned} \quad (19)$$

where $f_{self}$ is a multilayer perceptron, $G_{ea}$ is the same *erase*-followed-by-*add* mechanism used in DKVMN, and $G_{gru}$ is the gated recurrent unit (GRU) gate [62]. Moreover, $f_{neighbor}$ defines the information propagation to neighboring nodes based on the knowledge graph structure.

In the *predict* module, GKT predicts the student's performance at the next time step according to the updated temporal knowledge state:

$$y_k^t = \sigma(W_o h_k^{t+1} + b_k), \quad (20)$$

where $W_o$ is the weight parameter and $b_k$ is the bias term.

Recently, in an attempt to further explore knowledge structure, Tong et al. [63] proposed structure-based knowledge tracing (SKT), which aims to capture the multiple relations in knowledge structure to model the influence propagation among concepts. SKT was mainly motivated by an education theory, *transfer of knowledge* [64], which claims that students' knowledge state on some relevant KCs will also be changed when they are practicing on a specific KC due to the potential knowledge structure among KCs. Therefore, a student's knowledge state is determined by not only the temporal effect from the exercise sequence, but also the spatial effect from the knowledge structure. To concurrently model the latent spatial effects, SKT presented the synchronization and partial propagation methods to characterize the undirected and directed relations between

KCs, respectively. In this way, SKT could model influence propagation in the knowledge structure with both temporal and spatial relations. By introducing education theory, SKT can get more interpretable evolving knowledge state.

# 4 VARIANTS OF KNOWLEDGE TRACING MODELS

So far, we have presented all basic KT models in our taxonomy. Generally speaking, these basic models are usually based on simplified assumptions (e.g., the two-state assumption in BKT), and mainly leverage the learning interactions (i.e., exercises and responses) to estimate the knowledge state of students. However, students' real learning process is not simply represented by exercises and responses, but is influenced by many factors. As a consequence, the above basic KT models are straightforward, but have reduced performance in real-world learning scenarios. Many variants have been proposed to capture more complete learning process, which compensate for such performance loss to some extent. In the following, according to different learning phases, we classify the current variants of basic KT models into three categories and review them detaily: (1) modeling individualization before learning, (2) incorporating engagement and utilizing side information during learning , and (3) considering forgetting after learning.

## 4.1 Modeling Individualization before Learning

Everything and everyone has unique characteristics. For example, Liu et al. [65] considered several personalized factors (e.g., spatial and temporal preferences) of various tourists and proposed a "cocktail" approach to personalized travel package recommendation. Similarly, individualization in the KT task refers to that different students tend to have different learning characteristics (i.e., different learning rates or prior knowledge). However, most basic KT models fail to model individualization in this way. Yudelson et al. [66] noted that accounting for student-specific variability in the learning data could enhance the accuracy of BKT. In the following, we will introduce some variants that aim to consider individualization before learning in the KT task.

### 4.1.1 Modeling Individualization in BKT

The original BKT paper has discussed individualization. Specifically, for a specific KC, it takes advantage of the learning interactions of all students on it to learn its parameters. Similarly, for a specific student, all his/her learning interactions are utilized to fit his/her learning parameters [13]. In this way, BKT can learn different learning and performance parameters for different students and KCs. However, this approach achieves only a limited improvement compared with the original BKT model. When attempting to model individualization in BKT, the natural idea is to individualize the parameters in BKT for each student. Pardos and Heffernan [67] proposed two simple variants of BKT that individualize students' initial probability of mastery and the probability of transition from the unlearned state to the learned state, respectively. As shown in Fig. 10(a), a student node $S$ is added to individualize the initial probability of mastery $P(L_0[S])$ for each student. The student node has values that range from one up to the total number of students, which means that each student has their



(a) Bayesian knowledge tracing with individualized $P(L_0)$    (b) Bayesian knowledge tracing with individualized $P(T)$
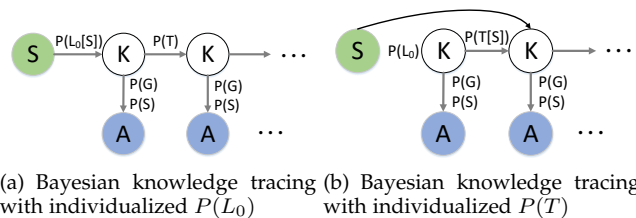
Fig. 10. The topology of Bayesian knowledge tracing with individualized initial probability or transition probability [67].

own initial probability of mastery. It also has a conditional probability table, which determines the probability that a specific student will have a specific node value. Similarly, if changing the connection of the student node from $P(L_0[S])$ to the subsequent knowledge nodes, the transition probability parameter can also be individualized. In this case, the student node gives individualized $P(T)$ parameters to each student, as shown in Fig. 10(b). Moreover, rather than individualizing only one kind of parameter in BKT, some other variants of BKT turn to individualize all four BKT parameters simultaneously[66]. Lee and Brunskill [68] indicated that when used in an intelligent tutoring system, the individualized BKT model can well improve the learning efficiency of students, which reduces about half amount of questions for $20\%$ of students to achieve mastery.

Another means of modeling individualization for a larger range of students is clustering, where we can train more appropriate models for different groups of students [69]. K-means is a basic clustering algorithm, which randomly initializes a set of K cluster centroids; these clusters are identified using Euclidean distance. Another popular clustering algorithm is spectral clustering, which represents the data as an undirected graph and analyzes the spectrum of the graph Laplacian obtained from the pairwise similarities of data points. Recently, some novel clustering algorithms have been proposed, including discrete nonnegative spectral clustering [70] and clustering uncertain data [71]. After clustering the students into K groups, we can train K different KT models and make predictions on the test data. The number of clusters K is then varied from K - 1 to 1 and the above process is repeated iteratively. Finally, we can obtain a set of K different predictions. Furthermore, for combining these predictions, there are two common methods [72]: (1) uniform averaging, which simply averages the K predictions, (2) weighted averaging, which combines the models by means of a weighted average.

### 4.1.2 Modeling Individualization in DKT

To model individualization in DKT, Minn et al. [73] proposed a model named deep knowledge tracing with dynamic student classification (DKT-DSC), which introduces individualization to DKT by exploiting the concept of clustering. According to students' previous performance, DKT-DSC assigns students with similar learning ability into the same group . The knowledge states of students in different groups are then traced by different DKT models. Moreover, considering the dynamic property of the learning ability, in order to dynamically assess the students' learning ability and accordingly assign students into new groups, each student's learning sequence is segmented into multiple time

intervals. At the start of each time interval, DKT-DSC will reassess students' learning ability and reassign their groups.

In more detail, students' learning ability is encoded as a vector with the length of the number of KCs that is updated after each time interval using all previous learning interactions on each KC. In DKT-DSC, the K-means clustering algorithm is utilized to split students with similar ability levels into the same group at each time interval. After learning the centroids of all K clusters, each student at each time interval $Seg_z$ is assigned to the nearest cluster by the following equation:

$$Cluster(Stu_i, Seg_z) = argmin \sum_{K}^{c=1} \sum_{d_{1:z}^i \in C_c} ||d_{1:z}^i - \mu_c||^2, \quad (21)$$

where $\mu_c$ is the mean of points in $C_c$ (a cluster set). The centroids for each cluster will remain unchanged throughout the entire clustering process. Through dynamic student clustering, DKT-DSC completes a more effective approach to realize the individualization in DKT.

In addition to modeling the individualization of students through extending existing KT models, Shen et al. [74] proposed a convolutional knowledge tracing model (CKT) to implicitly measure student individualization. More specifically, CKT considers two factors that influence students' individualization, that are individualized learning rates and individualized prior knowledge, respectively. Individualized learning rates represent that students have different absorptive capacities of knowledge. The sequence of student learning interactions can reflect different learning rates in the sense that students with high learning rates can rapidly master knowledge concepts, while others have to spend more time trying and failing. Therefore, it is reasonable to assess the differences in learning rate by simultaneously processing several continuous learning interactions within a sliding window of convolutional neural networks [24]. Moreover, individualized prior knowledge refers to students' prior mastery of knowledge, which is comprehensively assessed by students' historical learning interactions.

## 4.2 Incorporating Engagement during Learning

Student engagement is defined as *"the quality of effort students themselves devote to educationally purposeful activities that contribute directly to desired outcomes"* [75], which indicates a strong connection to students' learning process. Generally speaking, student engagement can impact or reflect student performance; in other words, higher engagement tends to result in better learning outcomes. Therefore, considering student engagement during learning process could potentially improve KT results [76]. In this section, we will present some variants that incorporate student engagement during learning into KT models.

### 4.2.1 Incorporating Engagement into BKT

Generally, student engagement is hard to measure directly, some online learning systems have made use of sensor data to measure student engagement. For example, inexpensive portable electroencephalography (EEG) devices can successfully help to detect a variety of student mental states related to the learning process, which can be seen as reflections of student engagement [77]. Xu et al. [78] proposed two types

of methods that combine EEG-measured mental states to improve the performance of BKT. Concretely, the first one inserts a one dimensional binary EEG measure into BKT, where the EEG-BKT structure extends BKT by adding a binary variable node $E$ between the knowledge node and the answer node. The second combined multi-dimensional continuous EEG measures in KT (EEG-LRKT), which uses logistic regression to combine an $m$-dimensional continuous variable $E$ extracted from the raw EEG signal in BKT.

However, in most cases, it is difficult to collect sensor data on every student. Therefore, Schultz and Arroyo [79] proposed the knowledge and affect tracing (KAT) model to parallelly model both knowledge and engagement. KAT is a sensorless model which does not rely on any sensor data. In this model, both knowledge and engagement are assumed to have direct influences on student performance. KAT takes three kinds of disengagement behaviors into account : quick guess (the student makes a attempt very quickly), bottom out hint (all available hints are used) and many attempts (making more than three attempts for an exercise). These three behaviors are grouped as "gaming" behaviors in order to predict students' knowledge and engagement at each learning interaction. Rather than assuming equal influence of knowledge and engagement on students' knowledge state, one variation on the KAT model defines the connection between knowledge and engagement, and accordingly considers that a student's knowledge state will influence their engagement. For example, students are more likely to disengage with knowledge they are not familiar. Moreover, rather than explicitly modeling student engagement, Schultz and Arroyo [80] further proposed the knowledge tracing with behavior (KTB) model, which has only one latent knowledge node as a combination of both knowledge and engagement. KTB assumes that both engagement and performance are expressions of a knowledge state. The Bayesian estimation of knowledge state needs to be infered by both student engagement and performance at every point in time.

### 4.2.2 Incorporating Engagement into DKT

To incorporate engagement into DKT, Mongkhonvanit et al. [81] proposed to add five features in the process of watching videos on the MOOC to the input of the DKT model. These features reflect student engagement from various aspects, that are playback speed, whether or not the video was paused, fast-forwarded or rewound, and whether or not the video was completed. For example, if a student watch a video in a much faster playback speed, he/she tends to impatient and absent-minded. Moreover, another two features: whether or not the exercise was submitted with an answer selected and whether or not the exercise was a part of an end-of-unit quiz, are also considered together. The experimental results indicate that DKT can achieve better performance simply through incorporating the above binarized engagement covariates.

## 4.3 Utilizing Side Information during Learning

Most KT models only leverage the performance data (i.e., exercises and student answers) to assess the students' knowledge state. Although these models have achieved

quite good results and have been successfully applied in online learning systems, there are still many other types of side information collected during learning process that can be utilized to obtain the students' knowledge state more precisely. In this part, we will introduce several variants that attempt to take advantage of these different types of side information during learning.

### 4.3.1 BKT with side information

Many models that apply side information in the BKT context have been proposed. In this section, we first introduce several works that extend BKT to enable modeling only one kind of side information for specific purposes, after which we present a general model that can utilize all types of side information.

Firstly, in terms of a student's first response time, a short initial response time could either indicate high proficiency or 'gaming' behavior, while a long first response time could either indicate careful thinking or lack of concentration. Since the connection between initial response time and knowledge state could be influenced by many complex and blended factors, Wang and Heffernan [82] proposed to discretize the continuous first response time into four categories (i.e., extremely short, short, long, extremely long) to eliminate unnecessary information and simplify the latent complex possibilities. It then builds a one-by-four parameter table for KT, in which each column represents the category of first response time of the previous exercise, while the relevant values represent the probability of correct answers.

Second, regarding tutor intervention, Beck et al. [83] proposed the Bayesian evaluation and assessment model, which simultaneously assesses student' knowledge states and evaluates the lasting impact of tutor intervention. More specifically, it adds one observable binary intervention node to BKT, where *True* means the tutor intervention occurs in corresponding interactions while *False* indicates the opposite. The connection between the intervention node and knowledge node indicates the potential impact of the tutor intervention on students' knowledge state. The intervention node is linked to all four BKT parameters. As a result, there are a total of eight parameters to learn in order to incorporate tutor intervention. One possible way of reducing the number of parameters is choosing to only link the intervention node to the learning rate parameter [84]. Similarly, Lin and Chi [85] developed the intervention-Bayesian knowledge tracing (Intervention-BKT) , which incorporates various types of instructional interventions into BKT and distinguishes their different effects on student performance. Specifically, the intervention node in the Intervention-BKT adds involves two types of interventions: *elicit and tell*. The relations between the intervention and performance nodes represent the impact of teaching interventions on a student's performance, while the relations between the intervention and knowledge nodes represent the impact of teaching interventions on a student's knowledge state. Therefore, at each learning interaction, while the student's knowledge state is conditional on both the previous knowledge state and the current intervention, the student's performance depends on both the present knowledge state and the current intervention.



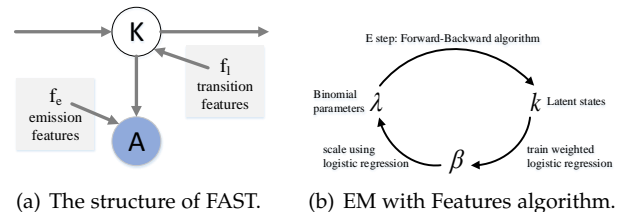(a) The structure of FAST.  (b) EM with Features algorithm.

Fig. 11. Feature Aware Student Knowledge Tracing Model [86].

Finally, rather than considering only one kind of side information, González-Brenes et al. [86] proposed a feature aware student knowledge tracing (FAST) model, an extension of BKT that allows to utilize all kinds of side information. Traditional BKT uses conditional probability tables for the guessing, slipping, transition and learning probabilities, meaning that the number of features involved in inference grows exponentially; therefore, as the number of features increases, the time and space complexity of the model also grow exponentially. To deal with this large number of features, FAST uses logistic regression parameters instead of conditional probability tables; thus, its number of features and complexity grow increase linearly rather than exponentially. Fig. 11(a) presents the graphical model structure of FAST, where $f_e$ are the features that parameterize the emission probabilities, while $f_l$ are the features that parameterize the transition probabilities. For parameter learning, FAST uses the Expectation Maximization with Features algorithm [87] and focuses on only emission features, as shown in 11(b). The E step uses the current parameter estimates $\lambda$ to infer the probability of the student having mastered the KC at each learning interaction. The parameters $\lambda$ are now a function of the weight $\beta$ and the feature vector $\boldsymbol{f}(t)$. $\boldsymbol{f}$ is the feature extraction function, while $\boldsymbol{f}(t)$ is the feature vector constructed from the observations at the $t-$th time step. The emission probability is represented with a logistic function, as follows:

$$\lambda(\beta)^{y',k'} = \frac{1}{1 + exp(-\beta^T \cdot \boldsymbol{f}(t))}, \qquad (22)$$

where $\beta$ is learned by training a weighted regularized logistic regression using a gradient-based search algorithm.

### 4.3.2 DKT with side information

With the goal of incorporating rich side information into DKT, Zhang et al. [88] proposed an extension to DKT that explored the inclusion of additional features to improve the DKT model. More specifically, it incorporates an auto-encoder network layer (a multi-layer neural network, as shown in Fig. 12(b)) to convert the higher-dimensional input data into smaller representative feature vectors, thereby reducing both the resource requirement and time needed for training. Students' response time, opportunity count, and first action are selected as incorporated side information and all input features are converted into a fixed-length input vector. Fig. 12 illustrates the process of feature incorporation. First, all input features are converted into categorical data and represented as a sparse vector by means of one-hot encodings. These encoded features are concatenated together to construct the higher-dimensional input vector as follows (as shown in Fig. 12(a)):
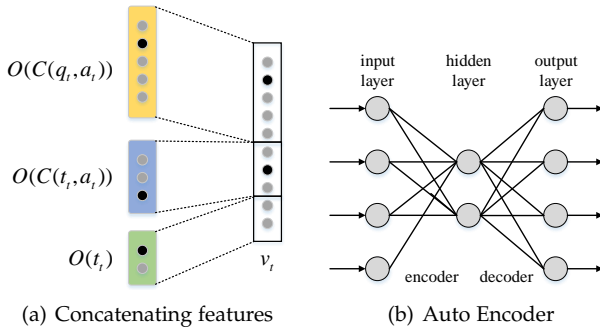
(a) Concatenating features      (b) Auto Encoder

Fig. 12. Feature processing for incorporating rich side information in to DKT [88].

$$C(e_t, a_t) = e_t + (max(e) + 1)a_t,$$
$$v_t = O(C(e_t, a_t)) \oplus O(C(t_t, a_t)) \oplus O(t_t), \quad (23)$$
$$v'_t = tanh(W_v v_t + b_v),$$

where $C()$ is the cross feature, $O()$ is the one-hot encoder format, $v_t$ represents the resulting input vector of each learning interaction, and $\oplus$ is the operation of concatenation. $W_v$ is the weight parameter and $b_v$ is the bias term. $e_t$ is the exercise, $a_t$ refers to the answer, and $t_t$ is the response time. Subsequently, auto-encoder is introduced to reduce the dimensionality without the loss of too much important information. Finally, the feature vectors extracted by auto-encoder will be the input of DKT.

To achieve more feasible integration of side information, Loh [89] proposed a deep knowledge tracing with decision trees (DKT-DT), which takes advantage of Classification And Regression Trees (CART) to preprocess the heterogeneous input features [90]. More specifically, CART is utilized to automatically partition the feature space and outputs whether or not a student can answer an exercise correctly. The predicted response and the true response are then encoded into a 4-bit binary code $O(f'_t, a_t)$; for example, $O(f'_t, a_t)$ is 1010 if the predicted response and the true response are both correct. $O(f'_t, a_t)$ is then concatenated with the original one-hot encoding of the exercise as the new input of DKT to train the corresponding model.

## 4.4 Considering Forgetting after Learning

In real-world scenarios, after learning, forgetting is inevitable [91]. The *Ebbinghaus forgetting curve theory* indicates that students' knowledge proficiency may decline due to forgetting [92]. Therefore, the assumption that a student's knowledge state will not change over time is not reasonable. Nevertheless, basic KT models, such as BKT, usually do not take forgetting into consideration. Recently, Huang et al. [93] proposed knowledge proficiency tracing (KPT) to model students' knowledge proficiency with both learning and forgetting theories, which can dynamically capture the change in students' proficiency level on KCs over time and track them in an effective and interpretable manner. In the following, we will introduce some variants that have attempted to consider forgetting after learning for more precise prediction results.

### 4.4.1 Considering Forgetting in BKT

Qiu et al. [94] have discovered that BKT would consistently over-predict a student's answers when a day or more had elapsed since her last response. The reason behind is that

BKT assumes that student performance is the same no matter how much time has passed. To consider how student performance declines with time, these authors proposed a BKT-Forget model, which hypothesized that students may forget as days go by. In the BKT-Forget model, a time node is added to specify which parameters should be affected by a new day and the new day node is fixed with a prior probability of 0.2. The parameter $forget\_n$ is introduced to represent the forgetting rate on a new day, while $forget\_s$ denotes the forget ting rate on the same day. However, although BKT-forget does consider the decline in student performance, it can only model forgetting occurred on a time scale of days. To model the continuous decay of knowledge as time progresses, Nedungadi and Remya [95] incorporated forgetting into BKT based on the assumption that learned knowledge decays exponentially over time [96]. An exponential decay function is thus utilized to update the knowledge mastery level. These authors further assumed that the chance of forgetting will increase if a student does not practice the knowledge concepts within 30 days. Moreover, Khajah et al. [41] introduced an approach that counts the number of intervening trials and treats each as an independent opportunity for forgetting to occur.

### 4.4.2 Considering Forgetting in PFA

Recall the PFA model in Eq.(5), in which the probability of students' mastery is estimated using a logistic function: $p(\theta) = \frac{1}{1+p(\theta)}$, where $\theta$ is defined as $\theta = \beta + \mu s + \nu f$. The original PFA model ignores the order of answers, not to mention the time between learning interactions. It is therefore difficult to directly incorporate time information into the original PFA model. Pelánek et al. [97] proposed PFAE (PFA Elo/Extended), a variant of the PFA model, which combines PFA with some aspects of the Elo rating system [98]. The Elo rating system was originally devised for chess rating (estimating players skills based on match results). In PFAE, $\theta$ is updated after each learning interaction, as follows:

$$\theta := \begin{cases} \theta + \mu \cdot (1 - p(\theta)) & \text{if the answer was correct,} \\ \theta + \nu \cdot p(\theta) & \text{if the answer was wrong.} \end{cases} \quad (24)$$

As the forgetting behavior of students is closely related to time, in order to consider forgetting, Pelánek [99] increased a time effect function $f$ to $\theta$, i.e., using $p(\theta + f(t))$ instead of $p(\theta)$, where $t$ is the time (in seconds) from the last learning interaction while $f$ is the time effect function.

### 4.4.3 Considering Forgetting in DKT

In order to represent the complex forgetting behavior, the DKT-forget model [100], which introduces forgetting into deep knowledge tracing, considers three types of side information related to forgetting: (1) the repeated time gap that represents the interval time between an interaction and the previous interaction with the same KC, (2) the sequence time gap that represents the interval time between the present interaction and the previous interaction in the sequence, and (3) past trial counts, which represent the number of times a student has attempted on exercises with the same KC. The two time gap features are used in minute scale, and all three features are discretized at $log_2$ scale. Those side information is concatenated as additional information

and represented as a multi-hot vector $c_t$, which is integrated with the embedding vector $v_t$ of the learning interaction, as follows:

$$v_t^c = \theta^{in}(v_t, c_t), \qquad (25)$$

where $\theta^{in}$ is the input integration function. The integrated input $v_t^c$ and the previous knowledge state $h_{t-1}$ is passed through the RNNs to update the student's knowledge state $h_t$ in the same way as in Eq.(7). The additional information at the next time step $c_{t+1}$ is also integrated with the updated student's knowledge state vector $h_t$:

$$h_t^c = \theta^{out}(h_t, c_{t+1}), \qquad (26)$$

where $\theta^{out}$ is the output integration function.

In a departure from considering forgetting in the framework of existing KT methods, Wang et al. [101] proposed a novel model, named HawkesKT, which introduced the Hawkes process to adaptively model temporal cross-effects in KT. The Hawkes process is good at modeling sequential events localized in time, as it controls corresponding temporal trends by the intensity function. The intensity function in HawkesKT is designed to characterize the accumulative effects of previous learning interactions, along with their evolutions over time. In HawkesKT, the temporal cross-effects and their evolution between historical learning interactions in a dynamic learning process.

## 5 Applications

Although knowledge tracing is an emerging research area, it has already been applied in a wide variety of scenarios. In the following, we survey the applications of KT models in three typical educational scenarios: learning resources recommendation, adaptive learning, and educational gaming.

### 5.1 Learning Resources Recommendation

Traditionally, learning resources for each student is selected in one of two ways. The first requires teachers to manually select suitable resources that matches students' knowledge levels for them to solve. However, this approach requires substantial time and efforts and different teachers may have different recommendations. The second allows students themselves to freely choose resources to learn; however, this may result in students choosing very easy or hard exercises that will not benefit their learning [104], leading to low learning efficiency. In recent years, the prevalence of intelligent tutoring systems and the development of KT methods make it possible to automatically recommend appropriate exercises to each student based on artificially designed intelligent algorithms.

Exercises are the most common learning resources in learning. Given the inferred knowledge levels of students, one common strategy is selecting the next exercise that will best advance the student's knowledge acquisition. Desmarais and Baker [104] proposed two extensions of the original BKT model, which considered exercises' difficulties or students' multiple-attempt behaviors respectively. These two extensions are integrated into a BKT-sequence algorithm to recommend exercises to students based on their knowledge state. More specifically, this BKT-sequence algorithm first determines the predicted range of scores for each exercise. For each exercise, it then computes an expected score that the student should get to achieve mastery, which is dependent on the their current knowledge state (for instance, lower knowledge state will result in higher expected scores). Finally, the algorithm returns the exercise with a predicted score that is closest to that of the expected score. Therefore, as the knowledge state of a particular KC grows, more difficult exercises will be recommended, as harder exercises are associated with a lower predictive score. Experimental results have shown that students using the BKT-sequence algorithm were enable to solve more difficult exercises, obtained higher performance and spent more time in the system than those students who used the traditional approach. Moreover, students also expressed that the BKT-sequence algorithm was more efficient.

In addition to exercises, there are also some other types of multi-modal learning resources, such as videos and figures. Machardy [105] opts to utilize an adaptation of BKT to improve student performance prediction by incorporating video observation. Experimental verification demonstrates the impact of both using and eschewing video data, as well as the rate of learning associated with a particular video. In this way, these authors developed a method to help people to evaluate the quality of video resources. Futhermore, they proposed the Template 1 Video model to incorporate video observations into BKT, which adds video activity as additional independent observation nodes to the BKT model. This model accordingly considers the probability that a given video resource will impart mastery on a student, and the transition probability is conditional only on the presence of either a video or an exercise. Thus, the quality of the video can be determined by its promotion of learning and this model can be leveraged as a tool to help to evaluate and recommend the video resources.

When recommending learning resources, existing solutions primarily aim to choose a simple strategy to assign non-mastered exercises to students. While reasonable, it is too broad to advance learning effectively. Huang et al. [106] proposed three more beneficial and specific objectives, which are *review and explore*, *smoothness of difficulty level* and *student engagement*, respectively. In more detail, *Review and explore* considers both enhancing students' non-mastered concepts with timely reviews and reserving certain opportunities to explore new knowledge; *Smoothness of difficulty level* means that the difficulty levels of several continuous exercises should vary within a small range as students gradually learn new knowledge; *Student engagement* indicates students' enthusiasm during learning, the recommended exercises should be in line with to their preferences. In order to support online intelligent education with the above three domain-specific objectives, the further presented a more reasonable multi-objective deep reinforcement learning (DRE) framework. In more detail, DRE presented three corresponding novel reward functions to capture and quantify the effects of the above three different objectives. This DRE framework is a unified platform to optimize multiple learning objectives, where more reasonable objectives also can be included in future. Experimental results show that DRE can effectively learn from the students' learning records to optimize multiple objectives and adaptively recommend suitable exercises to students.

TABLE 2
A summary of different types of knowledge tracing models, including their variants and applications.

| Models | Taxonomy | | | Main Technique | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Basic models | Variants | Applications | Bayesian networks | Logistic regression | RNNs/LSTMs | Memory networks | Attention mechanism | Graph neural networks | Others |
| BKT [13] | ✓ | | | ✓ | | | | | | |
| DBKT [31] | ✓ | | | ✓ | | | | | | |
| LFA [22] | ✓ | | | | ✓ | | | | | |
| PFA [23] | ✓ | | | | ✓ | | | | | |
| KTM [32] | ✓ | | | | ✓ | | | | | ✓ |
| DKT [26] | ✓ | | | | | ✓ | | | | |
| DKVMN [33] | ✓ | | | | | | ✓ | | | |
| EKT [34, 35] | ✓ | | | | | ✓ | ✓ | ✓ | | |
| PEBG [51] | ✓ | | | | | | | | | ✓ |
| SAKT [28] | ✓ | | | | | | | ✓ | | |
| AKT [29] | ✓ | | | | | | | ✓ | | ✓ |
| RKT [60] | ✓ | | | | | | | ✓ | | ✓ |
| GKT [27] | ✓ | | | | | | | | ✓ | |
| BKT-Forget [94] | | ✓ | | ✓ | | | | | | |
| Nedungadi and Remya [95] | | ✓ | | ✓ | | | | | | ✓ |
| PFAE [99] | | ✓ | | | ✓ | | | | | ✓ |
| DKT-Forget [100] | | ✓ | | | | ✓ | | | | |
| Pardos and Heffernan [67] | | ✓ | | ✓ | | | | | | |
| Lee and Brunskill [68] | | ✓ | | ✓ | | | | | | |
| Wang et al. [101] | | ✓ | | | | | | | | ✓ |
| Clustered-KT [69] | | ✓ | | ✓ | | | | | | ✓ |
| DKT-DSC [73] | | ✓ | | | | ✓ | | | | ✓ |
| CKT [74] | | ✓ | | | | | | | | ✓ |
| KAT [79] | | ✓ | | ✓ | | | | | | |
| KTB [80] | | ✓ | | ✓ | | | | | | |
| EEG-BKT [78] | | ✓ | | ✓ | | | | | | |
| Mongkhonvanit et al. [81] | | ✓ | | | | ✓ | | | | |
| Wang and Heffernan [82] | | ✓ | | ✓ | | | | | | |
| Beck et al. [83] | | ✓ | | ✓ | | | | | | |
| Intervention-BKT [85] | | ✓ | | ✓ | | | | | | |
| FAST [86] | | ✓ | | ✓ | | | | | | ✓ |
| Zhang et al. [88] | | ✓ | | | | ✓ | | | | |
| Pardos et al. [15] | | | ✓ | ✓ | | | | | | |
| Schodde et al. [103] | | | ✓ | ✓ | | | | | | |
| Liu et al. [20] | | | ✓ | | | ✓ | | | | ✓ |
| Desmarais and Baker [104] | | | ✓ | ✓ | | | | | | |
| Machardy [105] | | | ✓ | ✓ | | | | | | |

## 5.2 Adaptive Learning

Adaptive learning broadly refers to *"a learning process in which the content taught, or the way such content is presented, changes or "adapts" based on individual student responses, and which dynamically adjusts the level or types of instruction based on individual student abilities or preferences"* [107]. In contrast to learning resources recommendation, adaptive learning needs to design efficient learning schemes and dynamic learning paths to organize learning resources for students based on specific knowledge structures. However, the key problem associated with realizing adaptive learning is the same, that is dynamically measuring students' knowledge state, in which KT has been widely applied to solve this problem.

The first attempt made to apply KT to adaptive learning was the ACT Programming Tutor (APT) [13], where students were asked to write short programs and BKT was utilized to estimate their evolving knowledge state. This tutor can present an individualized sequence of exercises to each student based on these estimated knowledge states until the student has "mastered" each rule. Similarly, to provide adaptive language tutoring, Schodde et al. [103] presented an extension of BKT to predict possible effects that different tutoring actions may have on children's knowledge state, which will be utilized to guide the seletion of proper next tutoring actions. In the experiments, children's real responses on the training courses indicated that they indeed learned words in the process of human-robot interactions. Besides, further analysis indicated that the adaptive model significantly improved their performance by personalized robot tutoring, in contrast to a randomized training strategy.

In recent years, Massive Open Online Courses (MOOCs) have become an emerging modality of learning, especially in higher-education. Pardos et al. [15] have adapted BKT on the edX MOOC platform. The research object was a 14-week online course, which includes weekly video lectures and corresponding lecture problems. BKT was applied to study students' learning phenomena and enhance their learning on this course. More specifically, in order to better adapt BKT to the learning platform, the original BKT was modified in different respects. First, due to the lack of labeled KCs, elements of the course structure were utilized to project KCs to exercises; here, the problems would be directly seen as the KCs and its questions would be seen as the exercises belonging to the KC. Second, in order to capture a various amount of students' knowledge acquisition at each attempt, the modified model assigned different guess and slip parameters to different attempt counts. Third, to deal with the problem of multiple pathways in the system, which reflected that the impacts on learning may come from various resources, they considered the resource influence on learning as a credit/blame inference problem.

Generally speaking, students' cognitive structures in-

clude both students' knowledge level and the knowledge structure (e.g., *one digit addition* is the prerequisite knowledge of *two digit addition*) of learning items. Therefore, adaptive learning should keep consistency with both students' knowledge level and the latent knowledge structure. Nevertheless, existing methods for adaptive learning often separately focus on either the knowledge levels of students (i.e., with the help of specific KT model) or the knowledge structure of learning items. To fully exploit the cognitive structure for adaptive learning, Liu et al. [20] proposed a Cognitive Structure Enhanced framework for adaptive Learning, named CSEAL. CSEAL conceptualized adaptive learning as a Markov Decision Process, it first utilized DKT to trace the evolving knowledge states of students at each learning step. Subsequently, the authors designed a navigation algorithm based on the knowledge structure to ensure the logicality and reasonability of learning paths in adaptive learning, which also reduces the search space in the decision process. Finally, CSEAL utilized the actor-critic algorithm to dynamically determine what should be learned next. In this way, CSEAL can sequentially identify the suitable learning resources for different students.

### 5.3 Educational Gaming

The above two kind of applications are most common for knowledge tracing in traditional education field. Recently, knowledge tracing is also extended to be applied in more general scenarios, such as educational gaming. Long and Aleven [108] conducted a classroom experiment comparing a commercial game for equation solving, *DragonBox*, with a research-based intelligent tutoring system, *Lynnette*. The results indicated that students who used *DragonBox* enjoyed the experience more, while students who used *Lynnette* performed significantly better on the test. Therefore, it is possible to enable students to learn effectively and happily by designing suitable educational games on an online learning platform. In educational gaming, the paradigm of tracing students' knowledge state can also work for player modeling. Here, player modeling, which is the study of computational models of players in games, aims to capture human player characteristics and cognitive features [109]. For instance, Fisch et al. [110] revealed that children engage in cycles of increasingly sophisticated mathematical thinking over the course of playing an online game. Kantharaju et al. [111] presented an approach to trace player knowledge in a parallel programming educational game, which is capable of measuring the current players' real-time state over the different skills required to play an educational game based only on in-game player activities.

### 5.4 Summary

To summarize, we list the categories and techniques of all aforementioned knowledge tracing methods in Table 2, including the basic KT models, the variants of KT models and the applications of KT models. Moreover, to better help researchers and practitioners implement knowledge tracing solutions, as well as to facilitate future research in this domain, we have collected some popular public datasets, algorithms and relevant resources for knowledge tracing at the following URL: https://github.com/bigdata-ustc.

## 6 FUTURE RESEARCH DIRECTIONS

This survey has reviewed the abundant current developments in the field of knowledge tracing, including their variants and typical applications, as completely as possible. Nonetheless, as knowledge tracing is a young but promising research area, a large number of research problems remain that need to be urgently solved. In this section, we discuss several potential future research directions.

### 6.1 Knowledge Tracing with Interpretability

Given the interaction sequence of each student, knowledge tracing aims to monitor students' changing knowledge state during the learning process and thus accurately predict her performance on future exercises. As it is difficult to obtain the true knowledge state of students, the performance of KT models is usually indirectly evaluated with reference to prediction tasks: that is, the higher the prediction precision of students' responses on future exercises, the better the performance of the KT model. Therefore, interpretability is typically not the major focus of existing KT models, especially for those deep-learning based models with the end-to-end learning characteristic. However, interpretability is of significant importance in the domain of education; for example, students usually care more about why a specific item is recommended rather than which/what item is recommended. More attention should therefore be paid to improving the interpretability of KT models. To this end, some educational theories can be considered, such as the *Rasch model* used in AKT [29] and the *transfer of knowledge* used in SKT [63]. Moreover, we could consider incorporating knowledge tracing along with some static cognitive diagnosis models. For instance, a neural cognitive diagnosis framework was recently proposed by Wang et al. [112] with the goal of obtaining more accurate and interpretable diagnostic results. By incorporating educational theories or static cognitive diagnosis, we expect to more accurately measure students' knowledge state at each learning step. Consequently, both accurate and interpretable knowledge tracing results can be achieved.

### 6.2 Knowledge Tracing with Continuous Responses

Most KT models simplify the learning environments. For example, they assume that students' answers are only binary (i.e., either 0 or 1). The continuous value of students' answers (e.g., those on subjective exercises) are usually omitted. In real-world scenarios, students may answer different types of exercises with either discrete or continuous responses, with the latter accounting for a large proportion that cannot be ignored. Simple binarization of the continuous responses introduces inevitable systemic errors to the estimation of students' knowledge states. It is therefore necessary to develop KT models that can handle continuous responses and further measure students' knowledge states from both discrete and continuous responses. For example, in the domain of cognitive diagnosis, Liu et al. [113] proposed a fuzzy cognitive diagnosis framework for both objective and subjective scenes, which combines fuzzy set theory and educational hypotheses to model students' knowledge proficiency.

## 6.3 Knowledge Tracing with Student' Feedback

Learning records are passive reflections of students' knowledge proficiency. By contrast, student feedback provides us with their proactive understanding about their knowledge states, which in turn yields direct and real indicators of their learning situation. However, there are few KT models that take advantage of training data related to student' feedback, even though it can play an important role in fixing the KT results [114]. Wang et al. [115] have noted that feedback plays a positive role in learning, which may promote transfer and retention in learning from worked-out examples. Therefore, incorporating student feedback is a promising avenue that may yield better results.

## 6.4 Knowledge Tracing with Less Learning Data

Moreover, the learning of high-quality KT models inevitably requires a substantial amount of data to guarantee training stability. However, practical educational scenarios often suffer from the cold-start problem and the data isolation problem: e.g., students' learning data tends to be distributed across different schools and is also highly proprietary, so that it is difficult to gather the data for training [116]. Therefore, potential methods of combining concepts such as federated learning or active learning to train novel KT models are also a promising research direction.

## 6.5 Knowledge Tracing for General User Modeling

Generally speaking, user modeling refers to tools for characterizing users' behaviors (e.g., frequent locations), personal information (e.g., age, gender, and occupation) and latent features (e.g., interests and abilities), which facilitate the provision of targeted services for different users [117]. As a type of latent feature modeling, user ability modeling (including knowledge tracing) diagnoses the proficiency of users (not only individuals, but also groups of individuals, like user teams and companies) on specific skills/concepts. Therefore, in addition to education, knowledge tracing can be generally applied in a number of domains for user modeling, such as games, sports and recruitment.

## 7 CONCLUSIONS

In this survey, we conducted a comprehensive overview of knowledge tracing. More specifically, we first proposed a new taxonomy from the technical perspective, which split the basic KT models into three categories: (1) probabilistic models, (2) logistic models and (3) deep learning-based models. Based on this taxonomy, we reviewed existing KT models comprehensively. We then introduced several variants of KT models designed to model more complete learning process in three different phases: (1) modeling individualization before learning, (2) incorporating engagement and utilizing side information during learning , and (3) considering forgetting after learning. Subsequently, we summarized the applications of KT in three common educational scenarios: learning resources recommendation, adaptive learning, and educational gaming. Finally, we outlined some potential future directions for this young but promising research field. We hope that this extensive survey of knowledge tracing can serve as a basic framework for both researchers and practitioners in future research.

## REFERENCES

[1] Colette Chabbott and Francisco O Ramirez. Development and education. In *Handbook of the Sociology of Education*, pages 163–187. Springer, 2000.

[2] Marguerite Wotto. The future high education distance learning in canada, the united states, and france: Insights from before covid-19 secondary data analysis. *Journal of Educational Technology Systems*, 2020.

[3] UNESCO. and Irina Bokova. *Rethinking education: Towards a global common good?* UNESCO Publishing, 2015.

[4] Eleanour Snow and Marlene Kaplan. The new five-year federal strategic plan in stem education: What's in it for science? In *AGU Fall Meeting Abstracts*, volume 2018, pages ED43B–01, 2018.

[5] Lu Huang, Yihe Zhu, Li Chen, and Yonghe Zheng. Scientometric analysis for f0701: Applications and grants in 2018. *Studies in Science of Science*, 37(6):977–985, 2019.

[6] Anders D Olofsson, Göran Fransson, and J Ola Lindberg. A study of the use of digital technology and its conditions with a view to understanding what 'adequate digital competence'may mean in a national policy initiative. *Educational Studies*, pages 1–17, 2019.

[7] Ebba Ossiannilsson. Sustainability: Special issue "the futures of education in the global context: Sustainable distance education". *Sustainability*, 07 2020.

[8] Yue Suo, N. Miyata, H. Morikawa, T. Ishida, and Yuanchun Shi. Open smart classroom: Extensible and scalable learning system in smart space using web service technology. *IEEE Transactions on Knowledge & Data Engineering*, 21(6):814–828, 2009.

[9] Tuan Nguyen. The effectiveness of online learning: Beyond no significant difference and future horizons. *MERLOT Journal of Online Learning and Teaching*, 11(2):309–319, 2015.

[10] Emanuel and J. Ezekiel. Online education: Moocs taken by educated few. *Nature*, 503(7476):342, 2013.

[11] KURT VanLEHN. The relative effectiveness of human tutoring, intelligent tutoring systems, and other tutoring systems. *Educational Psychologist*, 46(4):197–221, 2011. doi: 10.1080/00461520.2011.611369.

[12] A. Cully and Y. Demiris. Online knowledge level tracking with data-driven student models and collaborative filtering. *IEEE Transactions on Knowledge and Data Engineering*, 32(10):2000–2013, 2020.

[13] Albert T. Corbett and John R. Anderson. Knowledge tracing: Modeling the acquisition of procedural knowledge. *User Modelling and User-Adapted Interaction*, 4(4):253–278, 1995.

[14] C. E. Dowling and C. Hockemeyer. Automata for the assessment of knowledge. *IEEE Transactions on Knowledge and Data Engineering*, 13(3):451–461, 2001.

[15] Zachary A Pardos, Yoav Bergner, Daniel T Seaton, and David E Pritchard. Adapting bayesian knowledge tracing to a massive open online course in edx. *Educational Data Mining*, 13:137–144, 2013.

[16] Chencheng Li, Pan Zhou, Li Xiong, Qian Wang, and Ting Wang. Differentially private distributed online learning. *IEEE Transactions on Knowledge and Data Engineering*, 30(8):1440–1453, 2018.

[17] B. C. M. Fung, K. Wang, and P. S. Yu. Anonymizing classification data for privacy preservation. *IEEE Transactions on Knowledge and Data Engineering*, 19(5):711–725, 2007.

[18] Wentao Wang, Guowei Xu, Wenbiao Ding, Gale Yan Huang, and Zitao Liu. Representation learning from limited educational data with crowdsourced labels. *IEEE Transactions on Knowledge and Data Engineering*, 2020.

[19] J. Cowan. A learner centered approach to online education. *British Journal of Educational Technology*, 44(6):E221–E222, 2013.

[20] Qi Liu, Shiwei Tong, Chuanren Liu, Hongke Zhao, Enhong Chen, Haiping Ma, and Shijin Wang. Exploiting cognitive structure for adaptive learning. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 627–635, 2019.

[21] James H. Block and Robert B. Burns. Mastery learning. *Review of Research in Education*, 4(1):3–49, 1976.

[22] Hao Cen, Kenneth Koedinger, and Brian Junker. Learning factors analysis–a general method for cognitive model evaluation and improvement. In *International Conference on Intelligent Tutoring Systems*, pages 164–175. Springer, 2006.

[23] Philip I Pavlik Jr, Hao Cen, and Kenneth R Koedinger. Performance factors analysis–a new alternative to knowledge tracing. *Online Submission*, 2009.

[24] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learn-

ing. *nature*, 521(7553):436, 2015.

[25] Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989.

[26] Chris Piech, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas J Guibas, and Jascha Sohl-Dickstein. Deep knowledge tracing. In *neural information processing systems*, pages 505–513, 2015.

[27] Hiromi Nakagawa, Yusuke Iwasawa, and Yutaka Matsuo. Graph-based knowledge tracing: Modeling student proficiency using graph neural network. In *2019 IEEE/WIC/ACM International Conference on Web Intelligence (WI)*, pages 156–163. IEEE, 2019.

[28] Shalini Pandey and George Karypis. A self-attentive model for knowledge tracing. *arXiv preprint arXiv:1907.06837*, 2019.

[29] 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*, page 2330–2339, New York, NY, USA, 2020.

[30] Leena Razzaq, Jozsef Patvarczki, Shane F. Almeida, Manasi Vartak, Mingyu Feng, Neil T. Heffernan, and Kenneth R. Koedinger. The assistment builder: Supporting the life cycle of tutoring system content creation. *IEEE Transactions on Learning Technologies*, 2 (2):157–166, 2009.

[31] Tanja Käser, Severin Klingler, Alexander G Schwing, and Markus Gross. Dynamic bayesian networks for student modeling. *IEEE Transactions on Learning Technologies*, 10(4):450–462, 2017.

[32] Jill-Jênn Vie and Hisashi Kashima. Knowledge tracing machines: Factorization machines for knowledge tracing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 750–757, 2019.

[33] Jiani Zhang, Xingjian Shi, Irwin King, and Dit-Yan Yeung. Dynamic key-value memory networks for knowledge tracing. In *Proceedings of the 26th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 2017. doi: 10.1145/3038912.3052580.

[34] Yu Su, Qingwen Liu, Qi Liu, Zhenya Huang, Yu Yin, Enhong Chen, Chris Ding, Si Wei, and Guoping Hu. Exercise-enhanced sequential modeling for student performance prediction. *The 32nd AAAI Conference on Artificial Intelligence*, pages 2435–2443, 2018.

[35] Qi Liu, Zhenya Huang, Yu Yin, Enhong Chen, Hui Xiong, Yu Su, and Guoping Hu. EKT: Exercise-aware knowledge tracing for student performance prediction. *IEEE Transactions on Knowledge and Data Engineering*, 33(1):100–115, 2021.

[36] Xiaoqing Huang, Qi Liu, Chao Wang, Haoyu Han, Jianhui Ma, Enhong Chen, Yu Su, and Shijin Wang. Constructing educational concept maps with multiple relationships from multi-source data. In *2019 IEEE International Conference on Data Mining (ICDM)*, pages 1108–1113, 2019.

[37] Radek Pelánek. Bayesian knowledge tracing, logistic models, and beyond: an overview of learner modeling techniques. *User Modeling and User-Adapted Interaction*, 27(3-5):313–350, 2017.

[38] Nguyen Thai-Nghe, Lucas Drumond, Tomáš Horváth, Artus Krohn-Grimberghe, Alexandros Nanopoulos, and Lars Schmidt-Thieme. Factorization techniques for predicting student performance. In *Educational recommender systems and technologies: Practices and challenges*, pages 129–153. IGI Global, 2012.

[39] Nguyen Thai-Nghe, Lucas Drumond, Tomás Horváth, and Lars Schmidt-Thieme. Using factorization machines for student modeling. In *International Conference on User Modeling, Adaptation, and Personalization*, 2012.

[40] S. Rendle. Factorization machines. In *2010 IEEE International Conference on Data Mining*, pages 995–1000, 2010.

[41] Mohammad Khajah, Robert V Lindsey, and Michael C Mozer. How deep is knowledge tracing? *arXiv preprint arXiv:1604.02416*, 2016.

[42] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[43] Chun-Kit Yeung and Dit-Yan Yeung. Addressing two problems in deep knowledge tracing via prediction-consistent regularization. *arXiv preprint arXiv:1806.02180*, 2018.

[44] Xiaolu Xiong, Siyuan Zhao, Eric G Van Inwegen, and Joseph E Beck. Going deeper with deep knowledge tracing. *Educational Data Mining*, 2016.

[45] Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou,

et al. Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626):471–476, 2016.

[46] Ghodai Abdelrahman and Qing Wang. Knowledge tracing with sequential key-value memory networks. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, jul 2019.

[47] Zhenya Huang, Qi Liu, Enhong Chen, Hongke Zhao, Mingyong Gao, Si Wei, Yu Su, and Guoping Hu. Question difficulty prediction for reading problems in standard tests. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2017.

[48] Qi Liu, Zai Huang, Zhenya Huang, Chuanren Liu, Enhong Chen, Yu Su, and Guoping Hu. Finding similar exercises in online education systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1821–1830, 2018.

[49] Yu Yin, Qi Liu, Zhenya Huang, Enhong Chen, Wei Tong, Shijin Wang, and Yu Su. Quesnet: A unified representation for heterogeneous test questions. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1328–1336, 2019.

[50] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.

[51] Yunfei Liu, Yang Yang, Xianyu Chen, Jian Shen, Haifeng Zhang, and Yong Yu. Improving knowledge tracing via pre-training question embeddings. In Christian Bessiere, editor, *IJCAI-20*, pages 1577–1583. International Joint Conferences on Artificial Intelligence Organization, 7 2020.

[52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

[53] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North*. Association for Computational Linguistics, 2019.

[54] Xizhe Wang, Xiaoyong Mei, Qionghao Huang, Zhongmei Han, and Changqin Huang. Fine-grained learning performance prediction via adaptive sparse self-attention networks. *Information Sciences*, 545:223–240, 2021.

[55] Jia Zhu, Weihao Yu, Zetao Zheng, Changqin Huang, and Gabriel Pui Cheong Fung. *Learning from Interpretable Analysis: Attention-Based Knowledge Tracing*. 2020.

[56] Youngduck Choi, Youngnam Lee, Junghyun Cho, Jineon Baek, Byungsoo Kim, Yeongmin Cha, Dongmin Shin, Chan Bae, and Jaewe Heo. Towards an appropriate query, key, and value computation for knowledge tracing. L@S '20, page 341–344, 2020.

[57] Dongmin Shin, Yugeun Shim, Hangyeol Yu, Seewoo Lee, Byungsoo Kim, and Youngduck Choi. Saint+: Integrating temporal features for ednet correctness prediction. In *LAK2021: 11th International Learning Analytics and Knowledge Conference*, apr 2021.

[58] Youngduck Choi, Youngnam Lee, Dongmin Shin, Junghyun Cho, Seoyon Park, Seewoo Lee, Jineon Baek, Chan Bae, Byungsoo Kim, and Jaewe Heo. EdNet: A large-scale hierarchical dataset in education, 2020.

[59] Frederique M Lord. *Applications of Item Response Theory to Practical Testing Problems*. LAWRENCE ERLBAUM ASSCCIAATES, 1980.

[60] Shalini Pandey and Jaideep Srivastava. RKT: Relation-aware self-attention for knowledge tracing. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. ACM, oct 2020.

[61] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S Yu. A comprehensive survey on graph neural networks. *arXiv preprint arXiv:1901.00596*, 2019.

[62] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *ArXiv:1406.1078*, 2014.

[63] Shiwei Tong, Qi Liu, Wei Huang, Zhenya Huang, Enhong Chen, Chuanren Liu, Haiping Ma, and Shijin Wang. Structure-based knowledge tracing: An influence propagation view. In *Proceedings of the The 19th IEEE International Conference on Data Mining*, 2020.

[64] Royer and M. James. Theories of the transfer of learning. *Educational Psychologist*, 14(1):53–69, 1979.

[65] Qi Liu, Yong Ge, Zhongmou Li, Enhong Chen, and Hui Xiong.

Personalized travel package recommendation. In *IEEE 11th International Conference on Data Mining*, pages 407–416. IEEE, 2011.

[66] Michael V. Yudelson, Kenneth R. Koedinger, and Geoffrey J. Gordon. Individualized bayesian knowledge tracing models. In *Lecture Notes in Computer Science*, pages 171–180. Springer Berlin Heidelberg, 2013.

[67] Zachary A. Pardos and Neil T. Heffernan. Modeling individualization in a bayesian networks implementation of knowledge tracing. In *User Modeling, Adaptation, and Personalization*, pages 255–266. Springer Berlin Heidelberg, 2010.

[68] Jung In Lee and Emma Brunskill. The impact on individualizing student models on necessary practice opportunities. *International Educational Data Mining Society*, 2012.

[69] Zachary A Pardos, Shubhendu Trivedi, Neil T Heffernan, and Gábor N Sárközy. Clustered knowledge tracing. In *International Conference on Intelligent Tutoring Systems*, pages 405–410, 2012.

[70] Y. Yang, F. Shen, Z. Huang, H. T. Shen, and X. Li. Discrete nonnegative spectral clustering. *IEEE Transactions on Knowledge and Data Engineering*, 29(9):1834–1845, 2017.

[71] B. Jiang, J. Pei, Y. Tao, and X. Lin. Clustering uncertain data based on probability distribution similarity. *IEEE Transactions on Knowledge and Data Engineering*, 25(4):751–763, 2013.

[72] Shubhendu Trivedi, Zachary A Pardos, and Neil T Heffernan. Clustering students to generate an ensemble to improve standard test score predictions. In *International conference on artificial intelligence in education*, pages 377–384. Springer, 2011.

[73] Sein Minn, Yi Yu, Michel C. Desmarais, Feida Zhu, and Jill-Jenn Vie. Deep knowledge tracing and dynamic student classification for knowledge tracing. In *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, nov 2018.

[74] Shuanghong Shen, Qi Liu, Enhong Chen, Han Wu, Zhenya Huang, Weihao Zhao, Yu Su, Haiping Ma, and Shijin Wang. Convolutional knowledge tracing: Modeling individualization in student learning process. *SIGIR '20: The 43rd International ACM SIGIR conference on research and development in Information Retrieval Virtual Event China July, 2020*, pages 1857–1860, 2020.

[75] Vicki Trowler. Student engagement literature review. *The higher education academy*, 11(1):1–15, 2010.

[76] Robert M Carini, George D Kuh, and Stephen P Klein. Student engagement and student learning: Testing the linkages. *Research in higher education*, 47(1):1–32, 2006.

[77] Chris Berka, Daniel J Levendowski, Michelle N Lumicao, Alan Yau, Gene Davis, Vladimir T Zivkovic, Richard E Olmstead, Patrice D Tremoulet, and Patrick L Craven. Eeg correlates of task engagement and mental workload in vigilance, learning, and memory tasks. *Aviation, space, and environmental medicine*, 78(5): B231–B244, 2007.

[78] Yanbo Xu, Kai-min Chang, Yueran Yuan, and Jack Mostow. Eeg helps knowledge tracing! In *Workshop on Utilizing EEG Input in Intelligent Tutoring Systems (ITS2014 WSEEG)*, page 43, 2014.

[79] Sarah Schultz and Ivon Arroyo. Tracing knowledge and engagement in parallel in an intelligent tutoring system. In *Educational Data Mining*, 2014.

[80] Sarah E Schultz and Ivon Arroyo. Modeling the interplay between knowledge and affective engagement in students. *International Journal of People-Oriented Programming (IJPOP)*, 3(2):56–74, 2014.

[81] Kritphong Mongkhonvanit, Klint Kanopka, and David Lang. Deep knowledge tracing and engagement with moocs. In *Proceedings of the 9th International Conference on Learning Analytics & Knowledge*, pages 340–342, 2019.

[82] Yutao Wang and Neil T Heffernan. Leveraging first response time into the knowledge tracing model. *International Educational Data Mining Society*, 2012.

[83] Joseph E Beck, Kai-min Chang, Jack Mostow, and Albert Corbett. Does help help? introducing the bayesian evaluation and assessment methodology. In *International Conference on Intelligent Tutoring Systems*, pages 383–394. Springer, 2008.

[84] Michael Sao Pedro, Ryan Baker, and Janice Gobert. Incorporating scaffolding and tutor context into bayesian knowledge tracing to predict inquiry skill acquisition. In *Educational Data Mining*, 2013.

[85] Chen Lin and Min Chi. Intervention-bkt: incorporating instructional interventions into bayesian knowledge tracing. In *International conference on intelligent tutoring systems*, pages 208–218. Springer, 2016.

[86] José González-Brenes, Yun Huang, and Peter Brusilovsky. General features in knowledge tracing to model multiple subskills,

temporal item response theory, and expert knowledge. In *The 7th International Conference on Educational Data Mining*, pages 84–91. University of Pittsburgh, 2014.

[87] Taylor Berg-Kirkpatrick, Alexandre Bouchard-Côté, John DeNero, and Dan Klein. Painless unsupervised learning with features. In *The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 582–590, 2010.

[88] Liang Zhang, Xiaolu Xiong, Siyuan Zhao, Anthony Botelho, and Neil T Heffernan. Incorporating rich features into deep knowledge tracing. In *Proceedings of the Fourth (2017) ACM Conference on Learning@ Scale*, pages 169–172, 2017.

[89] Wei-Yin Loh. Classification and regression trees. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 1(1):14–23, 2011.

[90] Lap Pong Cheung and Haiqin Yang. Heterogeneous features integration in deep knowledge tracing. In *International Conference on Neural Information Processing*, pages 653–662. Springer, 2017.

[91] Shaul Markovitch and Paul D Scott. The role of forgetting in learning. In *Machine Learning Proceedings 1988*, pages 459–465. Elsevier, 1988.

[92] Yuying Chen, Qi Liu, Zhenya Huang, Le Wu, Enhong Chen, Runze Wu, Yu Su, and Guoping Hu. Tracking knowledge proficiency of students with educational priors. *The 26th ACM International Conference on Information and Knowledge Management (CIKM'2017)*, pages 989–998, 2017.

[93] Zhenya Huang, Qi Liu, Yuying Chen, Le Wu, Keli Xiao, Enhong Chen, Haiping Ma, and Guoping Hu. Learning or forgetting? a dynamic approach for tracking the knowledge proficiency of students. *ACM Transactions on Information Systems (TOIS)*, 38(2): 1–33, 2020.

[94] Yumeng Qiu, Yingmei Qi, Hanyuan Lu, Zachary A Pardos, and Neil T Heffernan. Does time matter? modeling the effect of time with bayesian knowledge tracing. In *Educational Data Mining*, pages 139–148, 2011.

[95] Prema Nedungadi and MS Remya. Incorporating forgetting in the personalized, clustered, bayesian knowledge tracing (pc-bkt) model. In *2015 International Conference on cognitive computing and information processing (CCIP)*, pages 1–5. IEEE, 2015.

[96] Geoffrey R Loftus. Evaluating forgetting curves. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 11(2): 397, 1985.

[97] Radek Pelánek, Jan Papoušek, and Vít Stanislav. Adaptive practice of facts in domains with varied prior knowledge. In *Proceedings of the 7th International Conference on Educational Data Mining*, pages 6–13, London, United Kingdom, 2014. International Educational Data Mining Society. ISBN 978-0-9839525-4-1.

[98] Radek Pelánek. Application of time decay functions and elo system in student modeling. In *Educational Data Mining*, 2014.

[99] Radek Pelánek. Modeling students' memory for application in adaptive educational systems. *International Educational Data Mining Society*, 2015.

[100] Koki Nagatani, Qian Zhang, Masahiro Sato, Yan-Ying Chen, Francine Chen, and Tomoko Ohkuma. Augmenting knowledge tracing by considering forgetting behavior. In *The World Wide Web Conference*. ACM, may 2019.

[101] Chenyang Wang, Weizhi Ma, Min Zhang, Chuancheng Lv, Fengyuan Wan, Huijie Lin, Taoran Tang, Yiqun Liu, and Shaoping Ma. Temporal cross-effects in knowledge tracing. The International Conference on Web Search and Data Mining, page 517–525, New York, NY, USA, 2021. Association for Computing Machinery.

[102] Joseph E Beck and June Sison. Using knowledge tracing in a noisy environment to measure student reading proficiencies. 16 (2):129–143, 2006.

[103] Thorsten Schodde, Kirsten Bergmann, and Stefan Kopp. Adaptive robot language tutoring based on bayesian knowledge tracing and predictive decision-making. mar 2017.

[104] Michel C Desmarais and Ryan S Baker. A review of recent advances in learner and skill modeling in intelligent learning environments. *User Modeling and User-adapted Interaction*, 22(1): 9–38, 2012.

[105] Zachary Machardy. Applications of Bayesian Knowledge Tracing to the Curation of Educational Videos. *Electrical Engineering and Computer Sciences University of California at Berkeley Technical*, 2015.

[106] Zhenya Huang, Qi Liu, Chengxiang Zhai, Yu Yin, and Guoping Hu. Exploring multi-objective exercise recommendations in online education systems. In *The 28th ACM International Conference*

*on Information and Knowledge Management*, 2019.

[107] Steven Oxman and William Wong. White Paper: Adaptive Learning Systems. *DV X Innovations DeVry Education Group*, (February):1–30, 2014.

[108] Yanjin Long and Vincent Aleven. Educational game and intelligent tutoring system: A classroom study and comparative design analysis. *ACM Transactions on Computer-Human Interaction*, 24(3): 1–27, 2017.

[109] Georgios N Yannakakis, Pieter Spronck, Daniele Loiacono, and Elisabeth André. Player modeling. *Dagstuhl Follow Ups*, 2013.

[110] Shalom M Fisch, Richard Lesh, Elizabeth Motoki, Sandra Crespo, and Vincent F Melfi. Children's mathematical reasoning in online games: Can data mining reveal strategic thinking? *Child Development Perspectives*, 5(2):88–92, 2011.

[111] P. Kantharaju, K. Alderfer, J. Zhu, B. Char, B. Smith, and S Ontaón. Tracing player knowledge in a parallel programming educational game. 2019.

[112] Fei Wang, Qi Liu, Enhong Chen, Zhenya Huang, Yuying Chen, Yu Yin, Zai Huang, and Shijin Wang. Neural cognitive diagnosis for intelligent education systems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6153–6161, 2020.

[113] Qi Liu, Runze Wu, Enhong Chen, Guandong Xu, Yu Su, Zhigang Chen, and Guoping Hu. Fuzzy cognitive diagnosis for modelling examinee performance. *Acm Transactions on Intelligent Systems and Technology*, 9(4):1–26, 2018.

[114] Ekaterina Vasilyeva, Seppo Puuronen, Mykola Pechenizkiy, and Pekka Rasanen. Feedback adaptation in web-based learning systems. *International Journal of Continuing Engineering Education and Life Long Learning*, 17(4/5):337, 2007.

[115] Meng Wang, Zong Kai Yang, San Ya Liu, Hercy N H Cheng, and Zhi Liu. Using feedback to improve learning: Differentiating between correct and erroneous examples. In *International Symposium on Educational Technology*, 2015.

[116] Jinze Wu, Zhenya Huang, Qi Liu, Defu Lian, Hao Wang, Enhong Chen, Haiping Ma, and Shijin Wang. Federated deep knowledge tracing. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, WSDM '21, page 662–670, New York, NY, USA, 2021.

[117] G. Fischer. User modeling in human–computer interaction. *User Modeling and User-Adapted Interaction*, 11(1-2):65–86, 2001.