

# Beyond Black-Box Deep Knowledge Tracing: Transformers with Representational Grounding for Pedagogical Interpretability

Concha Labra <sup>1\*</sup> , Olga C. Santos <sup>2</sup>

<sup>1</sup> Universidad Nacional de Educación a Distancia: Madrid, Madrid, ES; clabra@dia.uned.es

<sup>2</sup> Universidad Nacional de Educación a Distancia: Madrid, Madrid, ES; ocsantos@dia.uned.es

\* Author to whom correspondence should be addressed.

## Abstract

This study introduces iDKT, an interpretable-by-design Transformer model that utilizes *Representational Grounding* to align deep latent representations with educational constructs, leveraging the high accuracy of deep knowledge tracing models while addressing their inherent lack of interpretability. We introduce a formal validation framework to verify the alignment of iDKT's internal representations and, using Bayesian Knowledge Tracing (BKT) as a reference, evaluate the model across multiple educational datasets. Results demonstrate that iDKT maintains state-of-the-art predictive performance while yielding additional interpretable insights at a significantly higher granularity than those provided by the reference model. Specifically, iDKT identifies student-level initial knowledge and learning velocities, providing mastery estimates that are more sensitive to the nuances of individual behavioral patterns than those produced by standard BKT. These individualized insights enable precise diagnostic placement and dynamic pacing, allowing adaptive learning environments to tailor instruction to each student's unique learning profile with enhanced precision. This work offers both a robust methodology for evaluating the interpretability of Transformer-based models and a practical tool for improving educational effectiveness through data-driven personalization.

**Keywords:** deep knowledge tracing; transformer; interpretability; Bayesian Knowledge Tracing; educational data analysis; personalized learning

## 1. Introduction

Knowledge Tracing [1] is a fundamental task in the fields of Artificial Intelligence in Education, Intelligent Tutoring Systems and Massive Open Online Courses. Its primary objective is to model a student's dynamic knowledge state over time based on their history of interactions with learning materials, enabling systems to predict future performance and provide personalized instruction. As educational environments become increasingly diverse and digital, the ability to accurately track and interpret student mastery has become a critical requirement for scalable, effective education.

Historically, the field has been dominated by two distinct paradigms. The first, exemplified by Bayesian Knowledge Tracing (BKT) and its variants [2], relies on probabilistic graphical models that explicitly represent knowledge states. BKT models are intrinsically interpretable, being based on parameters such as initial knowledge, learning rate, or slipping and guessing probabilities that map directly to pedagogical constructs, allowing educators to understand how they work and trust their decisions. However, this interpretability

Received:

Revised:

Accepted:

Published:

**Copyright:** © 2025 by the authors.

Submitted to *Journal Not Specified* for possible open access publication under the terms and conditions of the [Creative Commons Attribution \(CC BY\)](https://creativecommons.org/licenses/by/4.0/) license.

comes at the cost of a simplicity that often limits its predictive power, making them struggle to capture the complex, non-linear dependencies often present in educational datasets.

The second paradigm emerged with the advent of Deep Knowledge Tracing (DKT) [3], which uses different variants of deep learning techniques from the initial Recurrent Neural Networks to current Transformers [4] to model student interactions. These models have achieved state-of-the-art predictive performance, significantly outperforming classical approaches by leveraging the high capacity of deep learning models that allows them to learn complex patterns [5]. Yet, this predictive power has come at a significant cost: interpretability. Deep learning models are notoriously opaque "black boxes," where the learned representations are distributed across high-dimensional latent spaces that bear no direct correspondence to constructs with a clear semantic meaning. This lack of transparency creates a trust gap for practitioners, who cannot easily discern why a model predicts a student has failed or succeeded, nor can derive actionable pedagogical insights from the model's internal state [6].

Current efforts to bridge this gap typically rely on post-hoc explainability methods, such as weights visualization or perturbation analysis [7,8]. While valuable for debugging, these techniques often provide only a superficial view of the model's decision-making process and do not guarantee that the learned representations align with meaningful constructs. Moreover, their application and interpretation require technical deep learning expertise, limiting their accessibility to practitioners without this specialized knowledge.

To address these limitations, we propose a shift towards interpretability-by-design, inspired by the emerging paradigm of Theory-Guided Data Science (TGDS) [9]. In TGDS, maintaining consistency with theoretical postulates is an architectural constraint rather than an afterthought. By integrating extensive domain knowledge, TGDS-based models can be constrained to learn representations that are both theoretically plausible and highly predictive. While this approach has been applied mostly to science—and specifically to physics [10]—we adapt it here to the educational domain.

Standard TGDS implementations typically rely on auxiliary loss functions to incorporate formal knowledge expressed as rules, algebraic constraints, or differential equations [11]. We propose a novel approach called *Representational Grounding* that, in contrast, utilizes auxiliary losses operating on projections of the Transformer's embeddings. This mechanism enables the model to learn representations that are consistent with semantically meaningful constructs.

The major contributions of this work are as follows:

- Proposal of Representational Grounding, a novel method that overcomes the black-box nature of Transformers by providing interpretability-by-design.
- Introduction of a formal validation framework to quantify interpretability via representational alignment, enabling a systematic characterization of the trade-off between reference fidelity and predictive performance.
- Application of Representational Grounding to the development of iDKT, a new type of knowledge tracing models that leverage the high accuracy inherent in deep learning while achieving pedagogical interpretability.
- Empirical demonstration of iDKT benefits by showing how it captures granular, student-specific insights—such as individualized initial knowledge and learning rates—that are beyond the capabilities of simpler models such as BKT.

The remainder of this paper is structured as follows. Section 2 reviews the current state of deep knowledge tracing and interpretability. Section ?? describes the proposed iDKT architecture and the Representational Grounding framework. Section 8 presents the experimental validation and answers the research questions. Finally, Section 9 discusses the results and their implications.

## 2. Related Work

### 2.1. Deep Knowledge Tracing

### 2.2. Deep Learning Interpretability

### 2.3. Theory-Guided Data Science

### 2.4. Individualized Bayesian Knowledge Tracing

#### 2.4.1. Parameters of the Vanilla Model

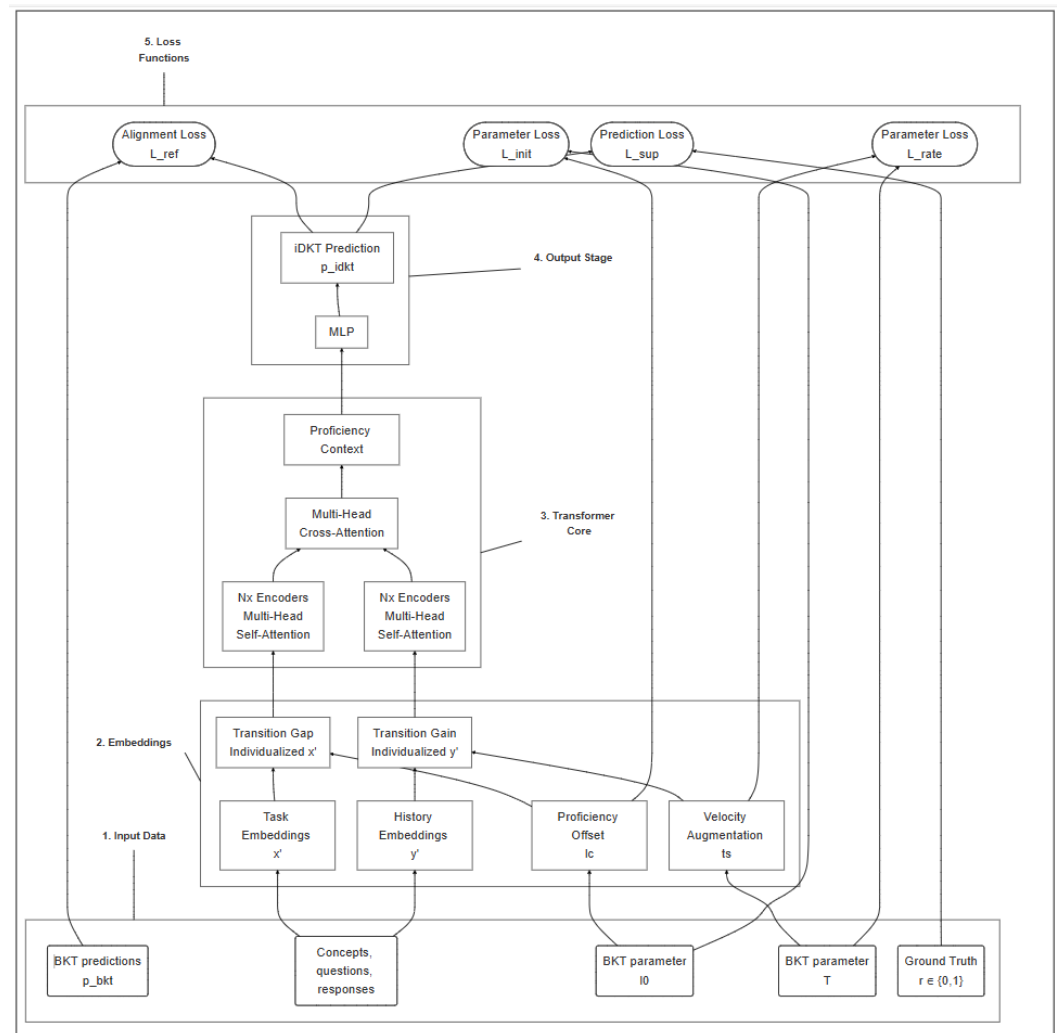
#### 2.4.2. Parameters Individualization

## 3. iDKT Model

We propose an Interpretable Deep Knowledge Tracing (iDKT) model with a Transformer-based architecture designed to bridge the gap between the high predictive capacity of deep learning and the intrinsic interpretability of simpler models such as Bayesian Knowledge Tracing (BKT). Unlike standard black-box models, iDKT utilizes a novel mechanism called Representational Grounding to anchor its latent representations to the conceptual space of an interpretable model chosen as reference.

The core architecture of iDKT, as illustrated in Figure 1, extends the standard Transformer framework [4] by incorporating specialized components for Representational Grounding, primarily integrated within the embedding layers and the multi-objective loss pipeline.

1. **Input Data:** This stage handles the ingestion of student interactions (concepts, questions, and binary responses  $r \in \{0, 1\}$ ). Crucially, it also loads values from the BKT reference model, including performance predictions ( $p_{BKT}$ ) and per-skill parameters such as initial mastery ( $l_0$ ) and learning transition rates ( $T$ ), which serve as grounding targets for the model.
2. **Embeddings:** Task and history information are projected into a continuous embedding space. We then apply the individualization transformations described in Section 4 where base embeddings are augmented by student-specific parameters: a Proficiency Offset ( $l_c$ ) and a Rate Augmentation ( $t_c$ ). This produces the Transition Gap (individualized task representation  $x'$ ) and the Transition Gain (individualized interaction history  $y'$ ), effectively defining the difficulty of a task relative to the student's baseline proficiency.
3. **Transformer Core:** The model employs a dual-encoder architecture followed by a cross-attention decoder. One encoder processes the individualized interaction history ( $y'_{1:t-1}$ ) to capture global student behavior, while a task encoder processes the current individualized task ( $x'_t$ ). The decoder uses multi-head cross-attention to synthesize these streams into a latent Proficiency Context, representing the student's current specialized knowledge state for the target task.
4. **Output Stage:** The latent Proficiency Context is passed through a multi-layer perceptron (MLP) that maps the deep representations to the final output space. The final output is the iDKT prediction ( $p_{iDKT}$ ), representing the probability that the student will respond correctly to the current task.
5. **Loss Functions (Grounding Pipeline):** During training, the architecture is supervised not only by the prediction loss ( $L_{sup}$ ) against ground truth outcomes but also by an alignment pipeline. This includes the Alignment Loss ( $L_{ref}$ ) which penalizes deviations from the BKT prediction, and Parameter Losses ( $L_{init}, L_{rate}$ ) that ground the model's internal proficiency and rate parameters to their BKT-derived theoretical counterparts.



**Figure 1.** The iDKT Architecture. The diagram illustrates the five functional stages: (1) Input Data ingestion including BKT targets, (2) Individualized Embeddings, (3) the Transformer Core, (4) the MLP-based Output, and (5) the Loss Functions for representational grounding.

## 4. Embeddings

In standard educational datasets, such as ASSISTments 2009, ASSISTments 2015, Algebra 2005, and others [12], student interactions are recorded at the level of specific questions or tasks, each of which is associated with one or more underlying concepts or knowledge components. This structure reflects the fact that proficiency in a concept (e.g., the Pythagorean Theorem) is acquired through interactions with a diverse range of tasks. While all tasks involving a concept share the same semantic core, they differ in their specific manifestations—most notably in their intrinsic difficulty or complexity. Therefore, an effective representation must capture both the shared identity of the concept and the unique deviation of the specific task.

In Transformer-based models [4] we can operationalize this principle representing the tasks as embedding vectors:

$$x' = c + u \cdot d \quad (\text{Task}) \quad (1)$$

In this formulation,  $c$  acts as the *Concept Anchor*, a vector representing the invariant semantic identity of the concept while the vector  $d$  represents the learnable *Question Variation Axis*, defining the specific direction of the "transition gap". The scalar  $u$  serves as the *Relational Magnitude*, representing the question's specific relative difficulty compared to other questions involving the same concept. Consequently, rather than encode arbitrary embeddings for every question, we encode the vector sum of two distinct components with clear semantic meaning: a base *concept identity* and a *difficulty shift*.

In a similar way, we can operationalize the interactions between questions and students as embedding vectors:

$$y' = e + u \cdot (f + d) \quad (\text{Interaction History}) \quad (2)$$

Here  $e$  represents the *Interaction Base*, which is a combined representation of concept  $c$  and the binary outcome  $r$  (correct/incorrect), while  $f$  represents the *Interaction Variation Axis*, which is similar to the Question Variation Axis ( $d_c$ ) but is specific to the interaction between a question and a student. The inclusion of  $d_c$  in the interaction shift ensures that the difficulty vectors are consistent across both questions ( $x'$ ) and interactions ( $y'$ ).

Extending this rationale, we can enrich the  $x'$  embeddings by integrating additional components with explicit semantic significance. Specifically, by adopting Bayesian Knowledge Tracing (BKT) as a reference model, we can incorporate vectors corresponding to its core theoretical parameters—Initial Knowledge ( $L_0$ ) and Learning Rate ( $T$ )—thereby grounding the deep representation in established pedagogical constructs.

To get individualized values for these parameters, we decompose them into population-level bases and student-specific deviations:

$$l_c = L_0 + k_s \cdot d_k \quad (\text{Personalized Initial Knowledge}) \quad (3)$$

$$t_c = T + v_s \cdot d_v \quad (\text{Personalized Learning Rate}) \quad (4)$$

where  $l_c$  is the personalized initial knowledge for concept,  $t_c$  is the personalized learning rate for concept,  $L_0$  and  $T$  are the population-level base embeddings,  $d_k$  and  $d_v$  are the learnable variation axes vectors (similar to the difficulty axis  $d$ ), and  $k_s$  and  $v_s$  are the scalar student-specific deviations learned for each individual.

We include these vectors to get the final input embedding for the encoder and decoder components of the Transformer:

$$x' = (c + u \cdot d) - l_c \quad (\text{Individualized Task}) \quad (5)$$

$$y' = (e + u \cdot (f + d)) + t_c \quad (\text{Individualized Interaction History}) \quad (6)$$

where  $c$  represents the concept embedding,  $u$  the question-specific difficulty shift,  $d$  the task variation axis,  $e$  the interaction base,  $f$  the interaction variation axis,  $l_c$  the personalized initial knowledge, and  $t_c$  the personalized learning rate.

The rationale for using difference for the individualized task ( $x'$ ) and sum for the interaction history ( $y'$ ) is due to their distinct semantic roles:

- $x'$  represents the *Transition Gap: Difficulty – Proficiency*. Under this relational logic, objective task difficulty is offset by prior proficiency, ensuring that task demands are defined relative to the subject's baseline. This formulation captures the difference between the task requirements and the current state, representing the residual gap after accounting for latent proficiency.
- $y'$  represents the *Transition Gain: Interaction + Rate*. The history encoder accumulates evidence from interactions, applying a consistent relational logic where the signal value is augmented by latent rate, ensuring that interaction outcomes are defined relative to the subject's pace. Under this formulation, the total value encompasses not only the interaction outcome but also the rate of progress through the state trajectory, as this incremental gain serves as a robust indicator of future performance. This approach, therefore, captures latent progression by augmenting observed outcomes with transition rate, thereby reflecting individualized acquisition rates.

## 5. Loss Functions

The model is trained using a multi-objective loss function designed to ensure that the high-capacity Transformer remains aligned with pedagogical principles through Representational Grounding. The total loss  $\mathcal{L}_{total}$  is defined as a weighted sum of different loss components described in detail below.

$$\mathcal{L}_{total} = L_{sup} + \lambda_{ref} L_{ref} + \lambda_{init} L_{init} + \lambda_{rate} L_{rate} + L_{reg} \quad (7)$$

### 5.1. Supervised Alignment ( $L_{sup}$ )

The primary objective  $L_{sup}$  uses standard Binary Cross-Entropy (BCE) between the iDKT performance predictions  $\hat{y}_t$  and the observed ground truth outcomes  $r_t \in \{0, 1\}$ . This loss ensures predictive accuracy by minimizing the deviance from observed student behavior:

$$L_{sup} = -\frac{1}{N} \sum_{t=1}^N [r_t \log(\hat{y}_t) + (1 - r_t) \log(1 - \hat{y}_t)] \quad (8)$$

### 5.2. Representational Grounding

The grounding losses ( $L_{ref}$ ,  $L_{init}$ ,  $L_{rate}$ ) use Mean Squared Error (MSE) to anchor deep representations to BKT-derived values. Specifically,  $L_{ref}$  forces behavioral predictions to stay close to the theoretical baseline, while  $L_{init}$  and  $L_{rate}$  ground the individualized parameters  $l_c$  and  $t_c$  in meaningful educational starting points and acquisition paces,

respectively. Instead of arbitrary latent weights, the model's internal states are projected through a sigmoid activation  $\sigma(\cdot)$  and compared directly to the reference values:

$$L_{ref} = \text{MSE}(\hat{y}, p_{BKT}) \quad (9)$$

$$L_{init} = \text{MSE}(\sigma(\bar{l}_c), L0_{BKT}) \quad (10)$$

$$L_{rate} = \text{MSE}(\sigma(\bar{t}_c), T_{BKT}) \quad (11)$$

where  $\bar{l}_c$  and  $\bar{t}_c$  are the average across the feature dimension of the individualized embeddings for proficiency and rate, respectively. This formulation forces the deep representation to be not only predictive but also semantically consistent with the reference constructs.

### 5.3. Inductive Bias Regularization ( $L_{reg}$ )

While the grounding losses anchor the global position of the latent space to the BKT parameter estimations,  $L_{reg}$  ensures that student-level individualization is *parsimonious*. This loss acts directly on the individualization parameters ( $u_q, k_s, v_s$ ) to ensure that the model only deviates from the theoretical prior when functionally necessary. We apply distinct  $L_2$  penalties to the scalar parameters governing variation:

$$L_{reg} = \lambda_u \sum_{q \in Q} u_q^2 + \lambda_k \sum_{s \in S} k_s^2 + \lambda_v \sum_{s \in S} v_s^2 \quad (12)$$

where  $u_q$  represents item difficulty,  $k_s$  is the student-specific knowledge gap, and  $v_s$  is the learning rate deviation. This formulation implements a *normal student prior*: the model assumes every subject adheres to the population-level parameters derived from the BKT reference unless their unique interaction history provides sufficient signal to justify the regularization cost.

## 6. Results

### 6.1. Experimental Setup

### 6.2. Datasets

### 6.3. Models

### 6.4. Research Questions

The experimental validation of iDKT is guided by the following research questions:

1. **Interpretability Validation (RQ1):** Is it possible to rigorously validate that a iDKT model, whose representations are grounded in a reference model, actually yields interpretable constructs?
2. **Accuracy–Interpretability Trade-Off (RQ2):** To what extent can deep knowledge tracing models be constrained for interpretability alignment without significantly degrading predictive accuracy?
3. **Higher Granularity (RQ3):** Can the proposed model capture individualized latent factors at a higher level of granularity than the reference model?

### 6.5. Interpretability Validation

To verify that iDKT's internal representations ( $l_c, t_c, u_q$ ) reflect educational constructs throughout training, we employ a verification framework based on two psychometric hypotheses:

1.  **$H_1$ : Convergent Validity (Latent Fidelity):** Pearson correlation between latent projections ( $l_c, t_c$ ) and reference BKT parameters. High alignment proves the model has internalized the theoretical constructs.



2.  **$H_2$ : Predictor Equivalence (Behavioral Alignment):** Functional substitutability of iDKT parameters into canonical BKT mastery recurrence equations. This ensures factors preserve their causal roles defined by theory.

#### 6.6. Accuracy–Interpretability Trade-Off

A key contribution of this work is the systematic exploration of the trade-off between predictive accuracy and theoretical fidelity. By modulating the grounding weight  $\lambda_{ref}$  in Equation 7, we identify the Pareto frontier of the model. This allows us to determine the "Inductive Bias Bonus"—points where theoretical grounding acts as a beneficial regularizer that improves generalization—and the point of "Over-Constraint," where excessive adherence to the reference model begins to degrade predictive power.

#### 6.7. Higher Granularity

To address RQ3, we introduce a metric termed *Individualization Volume*, defined as the variance ( $\sigma^2$ ) of the student-specific latent parameters ( $l_c, t_c$ ) across the population for each skill. In the reference BKT model, these parameters are fixed per skill, implying a population variance of zero ( $\sigma \approx 0$ ). Therefore, the ability of iDKT to capture granularity is verified by measuring the magnitude of dispersion in its learned parameter distributions. A significant non-zero variance indicates that the model has successfully identified valid student-specific traits—such as variable learning rates—that are obscured by the population-level averaging of the baseline model. We visualize this by plotting the standard deviation profiles for all knowledge components in the curriculum.

### 7. Plots

We utilize four key visualizations to support our findings:

1. **Confidence Heatmaps** (*per\_skill\_alignment\_predictions.png*): Support **RQ1**. A grid of students vs. skills showing regions of agreement (Green) and divergence (Red) between iDKT and BKT, validating that individualization is localized and grounded.
2. **Mastery Plot Mosaics** (*mosaic\_top\_students.png*): Supports **RQ3**. Longitudinal tracking of probability trajectories for individual students. Deviations from the monotonic BKT curve demonstrate the model's granular responsiveness to specific interaction histories.
3. **Pareto Frontier** (*idkt\_pareto\_frontier.png*): Supports **RQ2**. A scatter plot of Predictive Accuracy (AUC) vs. Theoretical Fidelity ( $r$ ) across varying  $\lambda$  weights, identifying the optimal trade-off point.
4. **Correlation Plots** (*per\_skill\_correlation\_predictions.png*): Supports **RQ1**. Bar charts of parameter correlation per skill, statistically quantifying the degree of semantic alignment ( $H_1$ ).

### 8. Results

To evaluate the proposed iDKT model, we conducted a series of experiments across multiple educational datasets. Table 1 summarizes the alignment metrics across different grounding strengths ( $\lambda$ ).

We observe consistent **Convergent Validity** ( $M_1$ ), with the correlation between the model's projected initial mastery ( $l_c$ ) and the theoretical prior ( $L_0$ ) remaining above 0.96 throughout the sweep (see Table 1).



**Table 1.** Construct Validity and Performance across the Grounding Spectrum.

Grounding Strength ( $\lambda$ )	Test AUC	$M_1$	$M_2$	$M_3$
0.00 (Baseline)	0.8317	0.9993	0.2652	-0.0325
<b>0.10</b>	<b>0.8322</b>	<b>0.9838</b>	<b>0.2949</b>	<b>-0.0330</b>
0.30	0.7984	0.9691	0.3192	-0.0330
0.50	0.7740	0.9884	0.2828	-0.0331

We observe consistent **Convergent Validity** ( $M_1$ ), with the correlation between the model’s projected initial mastery ( $l_c$ ) and the theoretical prior ( $L_0$ ) remaining above 0.96 throughout the sweep. This confirms that the Representational Grounding mechanism successfully anchors the deep latent space to the reference theory. Furthermore, the **Discriminant Validity** ( $M_3$ ) remains stable at  $r \approx -0.03$ , proving that the model successfully disentangles “Student Knowledge Gap” ( $k_c$ ) from “Student Learning Velocity” ( $v_s$ ) as distinct, non-redundant traits.

8.1. Pareto Curve Analysis

Our analysis reveals a non-linear trade-off between predictive accuracy and theoretical fidelity. Contrary to the common assumption that interpretability imposes a performance penalty, we identified an **“Inductive Bias Bonus”** at moderate grounding levels ( $\lambda \approx 0.10$ ).

As shown in Table 1, the model with  $\lambda = 0.10$  achieves a Test AUC of **0.8322**, slightly outperforming the unconstrained baseline (0.8317). This suggests that the BKT-based regularization acts as a beneficial inductive bias, preventing the Transformer from overfitting to noise in sparse interaction histories. However, excessive grounding ( $\lambda > 0.30$ ) leads to a sharp decline in predictive performance as the model becomes over-constrained by the simplicity of the reference theory.

8.2. Granularity of Individualization

While standard BKT assigns a fixed “Learning Rate” ( $T$ ) to all students for a given skill, iDKT captures a rich distribution of **Individualized Learning Velocities** ( $t_s$ ). Figure ?? illustrates this “Delta Distribution” ( $\Delta = t_s - T$ ). We observe a visible right-skewed variance, indicating that for many skills, the Deep Learning model identifies “fast-track” learning trajectories that classical population-level models underestimate. This granularity allows for **Precise Diagnostic Placement**, distinguishing between students who lack initial knowledge (*low*  $l_c$ ) versus those who suffer from slow acquisition pace (*low*  $t_s$ ).

8.3. Longitudinal Mastery Dynamics

The practical impact of these individualized parameters is evident in the mastery analysis. When simulating the mastery acquisition of “Fast” vs. “Slow” learners on the same sequence of correct responses:

- **Standard BKT** predicts identical mastery curves for both students.
- **iDKT** projects distinct trajectories, where “Fast” learners reach the 95% mastery threshold significantly earlier (fewer interactions) than “Slow” learners.

This “Informed Divergence” validates that iDKT does not merely mimic BKT labels but leverages its transformer core to dynamically adjust the **Velocity of Mastery** based on the student’s historical profile, enabling truly adaptive pacing in intelligent tutoring scenarios.

9. Discussion

Authors should discuss the results and how they can be interpreted from the perspective of previous studies and of the working hypotheses. The findings and their implications

should be discussed in the broadest context possible. Future research directions may also be highlighted.

## 10. Conclusions

This section is not mandatory, but can be added to the manuscript if the discussion is unusually long or complex.

## 11. Patents

This section is not mandatory, but may be added if there are patents resulting from the work reported in this manuscript.

**Author Contributions:** For research articles with several authors, a short paragraph specifying their individual contributions must be provided. The following statements should be used “Conceptualization, X.X. and Y.Y.; methodology, X.X.; software, X.X.; validation, X.X., Y.Y. and Z.Z.; formal analysis, X.X.; investigation, X.X.; resources, X.X.; data curation, X.X.; writing—original draft preparation, X.X.; writing—review and editing, X.X.; visualization, X.X.; supervision, X.X.; project administration, X.X.; funding acquisition, Y.Y. All authors have read and agreed to the published version of the manuscript.”, please turn to the [CRediT taxonomy](#) for the term explanation. Authorship must be limited to those who have contributed substantially to the work reported.

**Funding:** Please add: “This research received no external funding” or “This research was funded by NAME OF FUNDER grant number XXX.” and “The APC was funded by XXX”. Check carefully that the details given are accurate and use the standard spelling of funding agency names at <https://search.crossref.org/funding>, any errors may affect your future funding.

**Institutional Review Board Statement:** In this section, you should add the Institutional Review Board Statement and approval number, if relevant to your study. You might choose to exclude this statement if the study did not require ethical approval. Please note that the Editorial Office might ask you for further information. Please add “The study was conducted in accordance with the Declaration of Helsinki, and approved by the Institutional Review Board (or Ethics Committee) of NAME OF INSTITUTE (protocol code XXX and date of approval).” for studies involving humans. OR “The animal study protocol was approved by the Institutional Review Board (or Ethics Committee) of NAME OF INSTITUTE (protocol code XXX and date of approval).” for studies involving animals. OR “Ethical review and approval were waived for this study due to REASON (please provide a detailed justification).” OR “Not applicable” for studies not involving humans or animals.

**Informed Consent Statement:** Any research article describing a study involving humans should contain this statement. Please add “Informed consent was obtained from all subjects involved in the study.” OR “Patient consent was waived due to REASON (please provide a detailed justification).” OR “Not applicable” for studies not involving humans. You might also choose to exclude this statement if the study did not involve humans.

Written informed consent for publication must be obtained from participating patients who can be identified (including by the patients themselves). Please state “Written informed consent has been obtained from the patient(s) to publish this paper” if applicable.

**Data Availability Statement:** We encourage all authors of articles published in MDPI journals to share their research data. In this section, please provide details regarding where data supporting reported results can be found, including links to publicly archived datasets analyzed or generated during the study. Where no new data were created, or where data is unavailable due to privacy or ethical restrictions, a statement is still required. Suggested Data Availability Statements are available in section “MDPI Research Data Policies” at <https://www.mdpi.com/ethics>.

**Acknowledgments:** In this section you can acknowledge any support given which is not covered by the author contribution or funding sections. This may include administrative and technical support, or donations in kind (e.g., materials used for experiments). Where GenAI has been used for purposes such as generating text, data, or graphics, or for study design, data collection, analysis, or

interpretation of data, please add “During the preparation of this manuscript/study, the author(s) used [tool name, version information] for the purposes of [description of use]. The authors have reviewed and edited the output and take full responsibility for the content of this publication.”

**Conflicts of Interest:** Declare conflicts of interest or state “The authors declare no conflicts of interest.” Authors must identify and declare any personal circumstances or interest that may be perceived as inappropriately influencing the representation or interpretation of reported research results. Any role of the funders in the design of the study; in the collection, analyses or interpretation of data; in the writing of the manuscript; or in the decision to publish the results must be declared in this section. If there is no role, please state “The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results”.

Abbreviations

The following abbreviations are used in this manuscript:

- MDPI Multidisciplinary Digital Publishing Institute
- DOAJ Directory of open access journals
- TLA Three letter acronym
- LD Linear dichroism

Appendix A

Appendix A.1

The appendix is an optional section that can contain details and data supplemental to the main text—for example, explanations of experimental details that would disrupt the flow of the main text but nonetheless remain crucial to understanding and reproducing the research shown; figures of replicates for experiments of which representative data are shown in the main text can be added here if brief, or as Supplementary Data. Mathematical proofs of results not central to the paper can be added as an appendix.

Table A1. This is a table caption.

Title 1	Title 2	Title 3
Entry 1	Data	Data
Entry 2	Data	Data

Appendix B

All appendix sections must be cited in the main text. In the appendices, Figures, Tables, etc. should be labeled, starting with “A”—e.g., Figure A1, Figure A2, etc.

References

1. Corbett, A.T.; Anderson, J.R. Knowledge tracing: Modeling the acquisition of procedural knowledge. *User modeling and user-adapted interaction* **1994**, *4*, 253–278.

2. Šarić Grgić, I.; Grubišić, A.; Gašpar, A. Twenty-five years of Bayesian knowledge tracing: a systematic review, 2022.

3. Piech, C.; Bassen, J.; Huang, J.; Ganguli, S.; Sahami, M.; Guibas, L.J.; Sohl-Dickstein, J. Deep knowledge tracing. *Advances in neural information processing systems* **2015**, *28*.

4. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, Ł.; Polosukhin, I. Attention is all you need. *Advances in neural information processing systems* **2017**, *30*.

5. Abdelrahman, G.; Wang, Q.; Nunes, B. Knowledge tracing: A survey. *ACM Computing Surveys* **2023**, *55*, 1–37.

6. Bai, X.; et al. A Survey of Explainable Knowledge Tracing, 2024.

7. Fantozzi, M.; et al. The Explainability of Transformers - Current Status and Directions. *arXiv preprint arXiv:2401.09202* **2024**.

8. Di Marino, S.; et al. Ante-Hoc Methods for Interpretable Deep Models: A Survey, 2025.

9. Karpadne, A.; Atluri, G.; Faghmous, J.; Steinbach, M.; Banerjee, A.; Ganguly, A.; Shekhar, S.; Samatova, N.; Kumar, V. Theory-guided Data Science: A New Paradigm for Scientific Discovery from Data. *IEEE Transactions on Knowledge and Data Engineering* **2017**, *29*, 2318–2331. 411
10. Willard, J.; Jia, X.; Xu, S.; Steinbach, M.; Kumar, V. Integrating Physics-Based Modeling With Machine Learning: A Survey, 2022. 412
11. Von Rueden, L.; Mayer, S.; Beckh, K.; Georgiev, B.; Giesselbach, S.; Heese, R.; Kirsch, B.; Pfrommer, J.; Pick, A.; Ramamurthy, R.; et al. Informed Machine Learning – A Taxonomy and Survey of Integrating Knowledge into Learning Systems. *IEEE Transactions on Knowledge and Data Engineering* **2021**, p. 1–1. 413
12. Liu, Z.; Liu, Q.; Chen, J.; Huang, S.; Tang, J.; Luo, W. pyKT: a python library to benchmark deep learning based knowledge tracing models. *Advances in Neural Information Processing Systems* **2022**, *35*, 18542–18555. 414

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content. 415