# Causal-Counterfactual Frameworks for Adaptive Learning Path Optimization: Integrating TD-BKT, Multi-Armed Bandits, and Socratic Explanations

## 1. The Epistemological Crisis in Black-Box Adaptive Learning

The contemporary landscape of educational technology is defined by a paradox of precision and opacity. We have entered an era where deep learning architectures, specifically Recurrent Neural Networks (RNNs) and Transformer-based models like Deep Knowledge Tracing (DKT), can predict a student's future performance with unprecedented accuracy. By ingesting vast sequences of interaction data—every click, pause, and keystroke—these systems map the intricate topography of student behavior. However, this predictive power comes at a significant epistemological cost. While these models can accurately forecast *what* a student will get wrong, they remain largely silent on *why* the failure occurs, or more importantly, *how* a specific pedagogical intervention might alter that trajectory. They operate as black boxes, generating optimal learning paths that are mathematically rigorous yet pedagogically inscrutable.
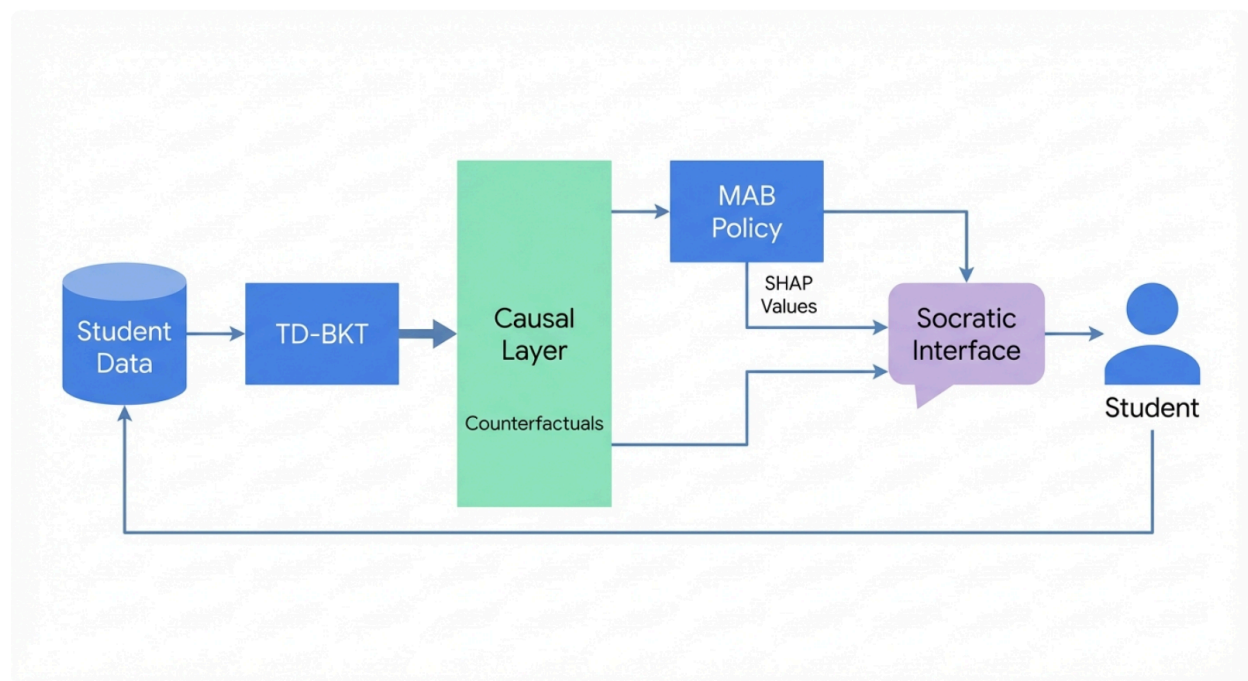
This lack of interpretability constitutes a critical barrier to the next generation of intelligent tutoring systems (ITS). In an educational context, a recommendation without an explanation is effectively a command, stripping the learner of agency. When a system directs a struggling student to "Review Module 4" without providing a causal rationale, it fosters dependency rather than metacognition. The student learns to follow the algorithm, not to understand their own cognitive gaps. Furthermore, standard predictive models are often correlational rather than causal. They might observe that students who watch a specific video tend to pass the exam, but they cannot definitively ascertain whether the video *caused* the success or if high-performing students simply tend to watch more videos. This distinction is not merely academic; it is the difference between a recommendation that works and one that wastes the learner's limited time.

To resolve this crisis, we must move beyond correlation and prediction into the realm of causality and counterfactual reasoning. We require systems capable of simulating alternative histories—answering the "What if?" questions that are central to human reasoning. "What if I had reviewed the prerequisites yesterday?" "What if I spent ten more minutes on this concept?" To achieve this, we propose a novel, hybrid architecture that integrates the probabilistic transparency of Time-Dependent Bayesian Knowledge Tracing (TD-BKT) with the

dynamic optimization of Multi-Armed Bandits (MAB), all governed by a Structural Causal Model (SCM).

This report details the theoretical and practical construction of such a system. It explores how we can formalize the learning process as a causal graph, allowing us to perform "algorithmic recourse"—identifying the specific, minimum-cost actions a student can take to reverse a negative outcome. We examine the use of SHAP (SHapley Additive exPlanations) to decompose the decisions of our optimization engine, providing transparency into the trade-offs between exploration and exploitation. Finally, we discuss the integration of these causal insights into a Socratic conversational agent, transforming raw data into a pedagogical dialogue that guides students toward self-regulated learning.

## Causal-Counterfactual Adaptive Learning Architecture



The architecture separates state estimation (TD-BKT) from policy optimization (MAB). The SCM layer intercepts the state estimates to run counterfactual simulations ('do-calculus'), feeding both the recommendation engine and the Socratic Agent. SHAP values are computed on the Bandit's arm selection policy to explain algorithmic choices.

# 2. Theoretical Foundations: Time-Dependent Bayesian Knowledge Tracing (TD-BKT)

The cornerstone of any counterfactual explanation system is a robust, interpretable model of

the learner's state. While neural approaches like Deep Knowledge Tracing (DKT) offer flexibility, their latent states are high-dimensional vectors that do not map one-to-one with pedagogical concepts. For a system designed to facilitate human understanding, the model's structure must mirror the cognitive structure of the domain. Bayesian Knowledge Tracing (BKT) provides this semantic clarity, and its Time-Dependent extension (TD-BKT) provides the necessary temporal dynamics to support actionable recourse.

## 2.1 The Limits of Classical BKT

Classical BKT models the learning of a skill as a Hidden Markov Model (HMM) with binary states: a skill is either mastered ($L_t = 1$) or not mastered ($L_t = 0$). The model updates its belief about the student's mastery based on observed binary responses (Correct/Incorrect) to questions tagged with that skill. The evolution of this state is governed by four fixed parameters per skill:

1. **Prior ($P(L_0)$):** The initial probability that the student knows the skill before the first practice opportunity.

2. **Learn ($P(T)$):** The probability that a student transitions from the unmastered state to the mastered state after a single practice opportunity.

3. **Guess ($P(G)$):** The probability that a student who has *not* mastered the skill will essentially guess correctly.

4. **Slip ($P(S)$):** The probability that a student who *has* mastered the skill will make a mistake.

While interpretable, classical BKT makes a strong assumption that is fatal for realistic counterfactuals: it assumes that learning events occur in a temporal vacuum. In standard BKT, the "Learn" probability is constant regardless of whether the practice opportunities happen five minutes apart or five weeks apart. It cannot account for the decay of memory over time, nor can it model the "spacing effect," where distributed practice yields better retention than massed practice. Consequently, a classical BKT model cannot answer questions like "What if I had studied this yesterday instead of today?" because "yesterday" and "today" are indistinguishable in the model's event-based time.

## 2.2 Incorporating the Physics of Memory: TD-BKT

Time-Dependent Bayesian Knowledge Tracing (TD-BKT) rectifies this by explicitly incorporating time as a variable in the causal graph. It draws upon the rich history of cognitive science, specifically the Ebbinghaus Forgetting Curve, which posits that memory retention decays exponentially over time in the absence of reinforcement.

In TD-BKT, the transition probabilities are no longer static scalars but dynamic functions of

the time interval $\Delta t$ since the last interaction.

The probability of a student *forgetting* a previously mastered skill is modeled as a function of time:

$$P(Forget|\Delta t) = 1 - e^{-\lambda_F \cdot \Delta t}$$

where $\lambda_F$ is a decay constant specific to the student or the skill complexity. Conversely, the probability of *learning* (transitioning from unmastered to mastered) can be modulated by the time between attempts, capturing the spacing effect. If $\Delta t$ is too short (cramming), the effective learning rate might be dampened; if optimal, it might be boosted.

This introduction of continuous time transforms the Bayesian network into a Dynamic Bayesian Network (DBN) where "Time" is an observed variable that causally influences the latent state transitions. From a structural modeling perspective, this is the breakthrough that enables temporal counterfactuals. By treating $\Delta t$ as a manipulable variable in our Structural Causal Model, we can simulate alternative timelines. We can hold the student's performance history constant but vary the timestamps of their interactions, recalculating the mastery probability for each hypothetical schedule. This allows the system to generate insights such as "Your retention of Geometry would be 12% higher if you had spaced these three review sessions over three days rather than completing them all in one hour."

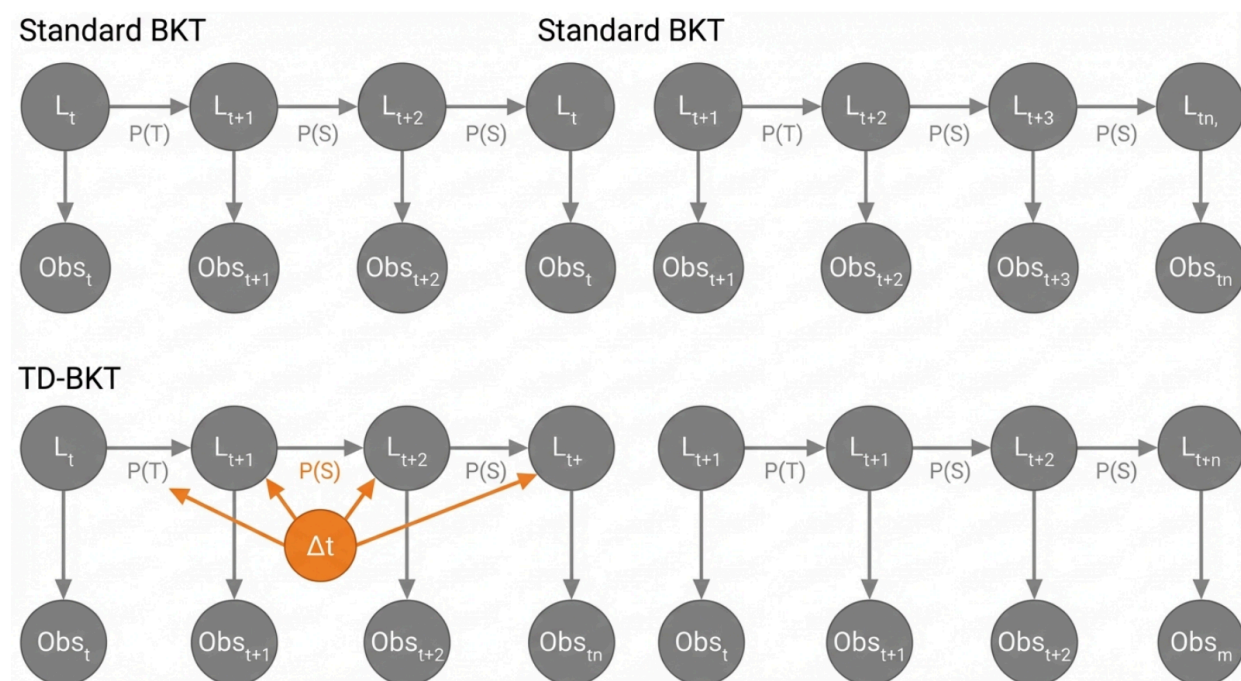## 2.3 Sensor Models for Complex Task Observation

The standard BKT assumption of binary "Correct/Incorrect" observations is increasingly insufficient for modern adaptive learning environments that include programming environments, virtual labs, or open-ended essays. In these contexts, a student's performance is multidimensional. They might solve a coding problem correctly but take three times longer than expected, or they might fail due to a simple syntax error (a slip) rather than a conceptual misunderstanding.

To handle this, our TD-BKT implementation integrates "sensor models"—auxiliary nodes in the Bayesian network that digest continuous and high-dimensional features. These sensors might track:

- **Response Time:** Logarithmic response times are strong indicators of fluency. Extremely fast responses on hard items often indicate guessing or "gaming the system."
- **Edit Distance:** In programming, how many edits did the student make? A high number of edits with a correct final result might suggest a "trial and error" strategy rather than true mastery.
- **Uncertainty/Hesitation:** Mouse tracking or pause analysis can reveal uncertainty even in correct answers.

These sensor nodes feed into the observation function $P(Obs|L)$, refining the estimate of the latent state. For counterfactual generation, this granularity is vital. It allows the system to distinguish between different types of failure. A failure driven by high "slip" probability (inferred perhaps from low response time) triggers a different counterfactual recommendation (e.g., "Slow down") than a failure driven by low mastery (e.g., "Review the concept"). This moves the feedback from a generic "You failed" to a precise diagnosis of the failure mode.

## Causal Graph Structure: BKT vs. TD-BKT



In Standard BKT (top), mastery transitions are purely event-driven. In TD-BKT (bottom), the time interval (Delta t) acts as a moderator variable, causally influencing both the probability of forgetting (decay) and the efficacy of learning (consolidation), enabling temporal counterfactuals.

# 3. Structural Causal Modeling (SCM) in Education

While TD-BKT provides a probabilistic map of the learner's state, it is essentially a descriptive model. To support true counterfactual reasoning—answering "What would have happened if...?"—we must formalize this probabilistic graph as a Structural Causal Model (SCM). The framework of SCM, pioneered by Judea Pearl, allows us to move beyond the first rung of the "Ladder of Causation" (Association) to the higher rungs of Intervention and Counterfactuals.

## 3.1 Formalizing the Learning SCM

We define our educational SCM as a tuple $\mathcal{M} = \langle U, V, F \rangle$.

- **Endogenous Variables ($V$):** These are the variables within the system's model that have explicit causal parents. In our context, $V$ includes:
  - $L_t$: The latent mastery state at time $t$.
  - $Obs_t$: The observed performance (correct/incorrect, response time).
  - $A_t$: The learning activity or intervention chosen (e.g., specific problem type, hint usage).
  - $\Delta t$: The time elapsed between interactions.

- **Exogenous Variables ($U$):** These represent the unobserved "background" factors that introduce stochasticity into the system. $U$ includes:
  - $U_{Student}$: Innate aptitude, current mood, or external distractions.
  - $U_{Content}$: Unmodeled variability in question difficulty or clarity.
  - $U_{Noise}$: Random measurement error in the sensors.

- **Structural Equations ($F$):** These functions describe the deterministic mechanisms by which the variables interact. For example, the mastery at time $t$ is determined by the previous mastery, the efficacy of the action taken, the time decay, and the student's specific learning noise:

$$L_t = f_L(L_{t-1}, A_{t-1}, \Delta t, U_{Student})$$

Crucially, in an SCM, the error term $U$ is not just "noise" to be averaged out; it is a specific, albeit unobserved, reality for that student at that moment.

## 3.2 The Logic of Counterfactuals: Abduction, Action, Prediction

The power of the SCM lies in its ability to compute counterfactuals through a precise three-step process. Let us consider a student who failed a quiz ($Obs_{fail}$). We want to ask: "Would they have passed if they had studied for 10 more minutes?"

1. **Abduction (The Diagnostic Step):** First, we condition the model on the observed reality ($Obs_{fail}$). We use the observed data to update the probability distribution of the exogenous variables $U$. We are essentially asking, "Given that the student failed, what is

the most likely state of their unobserved background factors (e.g., were they tired, or was the question essentially hard)?" This gives us a posterior distribution $P(U|Obs_{fail})$ specific to this event.

2. **Action (The Intervention Step):** Next, we perform an intervention on the model using the *do*-operator. We break the causal link that determined the original study time and force the variable to a new value: $do(\Delta t = \Delta t_{original} + 10 \text{ minutes})$. In the graph, this removes the arrows entering the $\Delta t$ node and fixes its value.

3. **Prediction (The Simulation Step):** Finally, we propagate this change through the network using the *inferred* background variables $U$ from step 1. We compute the new probability of mastery and success. Because we are using the specific $U$ inferred from the actual event, this is not a generic prediction for an average student; it is a specific prediction for *this* student, in *this* specific context, had they acted differently.

This process distinguishes a counterfactual from a simple conditional prediction. A conditional prediction ($P(Pass|Study + 10)$) asks, "What happens to the average student who studies 10 minutes more?" A counterfactual ($P(Pass_{Study+10}|Obs_{fail})$) asks, "Given that I *already* failed, what would have happened if I had studied 10 minutes more?" The latter is far more persuasive and relevant to the learner.

## 3.3 Causal Discovery and Graph Validation

A critical challenge in SCMs is defining the graph structure itself. How do we know that "Time Spent" causes "Mastery" and not the other way around (reverse causality, where mastery leads to faster completion)? In educational domains, we often rely on expert knowledge to define the initial graph skeletons—we know, biologically and pedagogically, that practice generally causes learning.

However, purely expert-driven graphs can miss subtle dependencies. To address this, we employ causal discovery algorithms applicable to time-series data, such as **PCMCI** (Peter-Clark Momentary Conditional Independence). PCMCI is designed to identify causal links in time-series data while controlling for autocorrelation, which is rampant in learning data (a student's performance today is highly correlated with their performance yesterday). By running PCMCI on aggregate student data, we can validate our expert-defined edges and potentially discover new ones—for instance, finding that "Hint Usage" in one session causally reduces "Response Time" in the next session (a positive transfer effect) or increases it (a dependency effect). This data-driven validation ensures that the counterfactuals generated by the system are grounded in empirical reality, not just theoretical assumptions.

# 4. Policy Optimization: The Multi-Armed Bandit Engine

While the TD-BKT model estimates the student's state and the SCM enables "what-if" reasoning, neither tells the system *what to do next* to maximize learning. This is an optimization problem. In our architecture, the decision-making engine is a Contextual Multi-Armed Bandit (CMAB). The "arms" of the bandit are the available learning activities (e.g., video, quiz, worked example), and the "context" is the student's current mastery state derived from the TD-BKT model.

## 4.1 Contextual Bandits and the Exploration-Exploitation Dilemma

In a standard A/B test, we might randomly assign students to different learning paths to see which is best on average. A Multi-Armed Bandit is a smarter, adaptive version of this. It constantly balances two competing goals:

- **Exploitation:** Selecting the educational activity that the model currently believes is the most effective for the student (maximizing immediate learning gain).
- **Exploration:** Selecting an activity about which the model is uncertain, to gather more data and potentially discover a better strategy (maximizing future system knowledge).

For personalized learning, we use **Contextual Bandits**. The policy $\pi(A_t|C_t)$ selects an action $A_t$ based on the context vector $C_t$, which includes the TD-BKT mastery probabilities, time since last review, and recent engagement metrics. The reward function $R_t$ is typically defined as the change in the student's total mastery across all knowledge components:

$$R_t = \sum_{k \in KCs} (P(L_{t+1}^k|Obs_{t+1}) - P(L_t^k|Obs_t))$$

## 4.2 Algorithm Selection: UCB vs. Thompson Sampling

To handle the exploration-exploitation trade-off, we employ the **Upper Confidence Bound (UCB)** algorithm. UCB calculates an "optimism bonus" for each action. For an arm $a$ and context $c$, the system estimates the expected reward $\hat{\mu}_{a,c}$ and adds a confidence interval term proportional to the uncertainty $\sigma_{a,c}$.

$$Score(a) = \hat{\mu}_{a,c} + \alpha \cdot \sigma_{a,c}$$

The system selects the arm with the highest score. If an arm has been tried few times (high uncertainty $\sigma$), its score is boosted, encouraging the system to "explore" it. As data accumulates, $\sigma$ shrinks, and the system naturally shifts toward "exploiting" the known best

strategy.

Alternatively, **Thompson Sampling** offers a Bayesian approach. Instead of calculating a distinct bound, it samples a reward probability from the posterior distribution of each arm. This probability matching strategy is often more robust in complex environments where rewards are delayed or noisy. In our architecture, we favor Thompson Sampling when the reward signal (learning gain) is noisy, as it handles the probabilistic nature of the BKT updates more gracefully than the deterministic bounds of UCB.

## 4.3 The "Abandonment" Constraint: MAB-A Models

A unique challenge in educational bandits is the problem of **abandonment**. In a standard web recommendation bandit (like for news articles), a bad recommendation just means a lost click. In education, a sequence of "optimal" but frustratingly difficult problems can cause a student to drop out entirely. A purely learning-maximizing bandit might relentlessly drill a student on their weakest area, causing burnout.

To mitigate this, we adopt a **Multi-Armed Bandit with Abandonment (MAB-A)** framework. The reward function is modified to include a penalty for disengagement. We model the probability of abandonment $P(Quit|A_t, C_t)$ as a separate outcome variable. The objective function becomes:

$$Maximize \sum_{t=0}^{T} \gamma^t \cdot (LearningGain_t - \lambda \cdot P(Quit_t))$$

where $\lambda$ is a hyperparameter balancing learning rigor against retention. This ensures that the generated learning paths are sustainable. When the Socratic agent explains a recommendation, it can implicitly reference this balance: "I'm suggesting a video review now—not because it's the fastest way to learn, but because you've been working hard on problems and my data suggests a break from active recall will help you sustain your focus."

# 5. Explaining the "Black Box": SHAP Attribution

Once the Bandit selects an activity, the system must explain *why* that specific arm was chosen over others. This is distinct from explaining the student's mastery; this is explaining the algorithm's choice. Since modern contextual bandits (especially those using Neural Networks or Deep RL) can be complex non-linear functions, we use **SHAP (SHapley Additive exPlanations)** to attribute the decision to specific features of the student's context.

## 5.1 SHAP for Bandit Arm Selection

SHAP values provide a game-theoretic distribution of credit for a model's output. In our context, the "game" is the selection of the optimal arm, and the "players" are the features of

the student's state (e.g., Mastery Level, Time Since Last Login, Quiz Score). The SHAP value $\phi_i$ represents how much feature $i$ pushed the estimated reward of the chosen arm away from the baseline average.

For example, if the Bandit recommends "Review Geometry," the SHAP analysis might reveal:

- $\phi_{GeometryMastery} = -0.4$ : Low mastery significantly increased the value of reviewing.
- $\phi_{ExamProximity} = +0.2$ : The fact that an exam is 2 days away added urgency.
- $\phi_{RecentFatigue} = -0.1$ : High recent activity slightly penalized the choice (risk of burnout), but was outweighed by the mastery need.

This decomposition allows the system to construct a narrative rationale: "I recommended this geometry review primarily because your mastery is below the threshold (-0.4 contribution) and your exam is approaching (+0.2 contribution)."
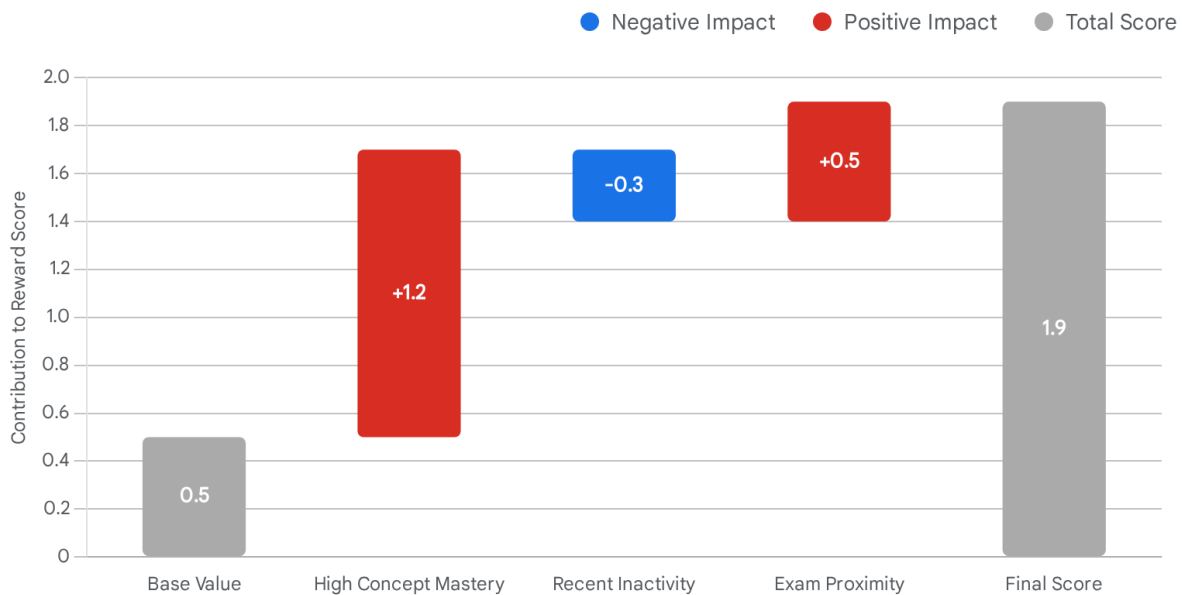
## 5.2 Addressing Sequential Dependencies with OrdShap

Standard SHAP assumes that features are independent, but in a learning path, the *order* of events is critical. Reviewing Algebra *before* Calculus is helpful; reviewing it *after* failing Calculus is a remedial action. The sequence matters. Standard SHAP might just say "Algebra Review" was important, without capturing the temporal interaction.

To address this, we utilize **OrdShap** (Ordered Shapley Values), a variation designed for sequential data. OrdShap disentangles the effect of a feature's value from its position in the sequence. It can distinguish between "You are struggling because you missed the quiz" (value effect) and "You are struggling because you took the quiz before watching the lecture" (position/sequence effect).

This capability is vital for generating "reordering" counterfactuals. If the OrdShap analysis shows a high negative contribution from the *position* of an activity, the system can generate a specific recommendation: "The data suggests that taking the quiz immediately after the lecture is less effective for you. If you had waited 24 hours (spacing effect), your predicted retention would be higher." This level of nuance—critiquing the *schedule* rather than just the *content*—is only possible with sequence-aware attribution.

# Why This Recommendation? SHAP Feature Attribution



The chart explains the Bandit's decision to recommend 'Advanced Integration.' Starting from the baseline expected reward, 'High Concept Mastery' (+1.2) and 'Exam Proximity' (+0.5) increase the score, while 'Recent Inactivity' (-0.3) slightly decreases it, resulting in a final high confidence score.

Data sources: SHAP Documentation, ArXiv, MDPI

# 6. Algorithmic Recourse: From Explanation to Action

Explanation explains the past; recourse changes the future. The ultimate goal of our system is not just to analyze learning but to optimize it through **Algorithmic Recourse**. Recourse refers to the ability of a person affected by a model's decision to change their situation to achieve a better outcome. In our context, it means giving the student a concrete plan to move from a "predicted failure" state to a "predicted success" state.

## 6.1 The Optimization Problem: Validity, Proximity, and Feasibility

We define the search for a counterfactual explanation as a constrained optimization problem.

Let $x$ be the student's current state vector (mastery, history, time spent) and $f(x)$ be the predictive model (TD-BKT) outputting the probability of success. We seek a counterfactual state $x'$ that minimizes a cost function:

$$Minimize \quad dist(x, x')$$

Subject to:

1. **Validity:** $f(x') \geq \tau$ (The new state must result in a passing probability above threshold $\tau$ ).

2. **Feasibility:** $x' \in \mathcal{F}$ (The changes must be possible). We cannot ask a student to change their past quiz scores or their innate aptitude. We *can* ask them to change their future study time or the order of upcoming tasks.

## 6.2 Effort-Aware Recourse

A naive distance metric (like Euclidean distance) might suggest that the "shortest path" to success is to increase mastery of a complex topic by 50% in one hour. Mathematically, this is a small shift in vector space; pedagogically, it is impossible.

To address this, we employ **Effort-Aware Distance Metrics**. We weight the features in the cost function based on the difficulty of changing them.

- **Immutable Features:** Past grades, demographics (Weight = $\infty$ ).
- **High-Effort Features:** Mastery of deep conceptual knowledge (Weight = High).
- **Low-Effort/Actionable Features:** Time spent on a task, utilizing a hint, reviewing a summary (Weight = Low).

The system prioritizes actionable features. It looks for the "low hanging fruit"—small behavioral changes that yield disproportionate gains in predicted mastery. This aligns with the concept of **Expected Minimum Cost (EMC)** recourse, where we seek a set of options that minimizes the expected effort the user must expend to achieve the goal.

## 6.3 Retrospective vs. Prospective Recourse

We categorize recourse into two distinct temporal modes, each serving a different pedagogical function.

**Retrospective Recourse ("What went wrong?")**

This analysis looks backward to identify missed opportunities. It is crucial for helping students understand the consequences of their study habits.

- *Example:* "If you had reviewed the 'Vector' module before attempting the 'Physics' quiz, your predicted score would have been 15% higher."
- *Mechanism:* The system uses the SCM to re-run the past sequence with alternative actions (interventions). It identifies critical decision points where a different choice would have significantly altered the outcome. This fosters reflection and helps students recognize the value of prerequisites.

**Prospective Recourse ("What to do next?")**

This analysis looks forward. It is a planning tool.

- *Example:* "Spending 10 more minutes on *Quadratic Equations* effectively guarantees (95% CI) a passing grade on tomorrow's final."
- *Mechanism:* The system projects the current state forward. It simulates different future study allocations (e.g., "Spend 10 mins here," vs "Spend 10 mins there"). It calculates the gradient of the success probability with respect to time for each topic. It then recommends the action with the steepest gradient—the most efficient use of the student's limited time.

For instance, consider a simulation where a student's mastery is tracked over five days. The actual trajectory (Baseline) might show a steady decline in retention due to lack of reinforcement, leading to a predicted failure on Day 5. However, the counterfactual simulation, triggered by a hypothetical intervention of 'Review Algebra' on Day 3, reveals a sharp divergence. The simulation might show that this single intervention would have prevented the mastery drop observed on Day 5, leading to a final mastery score 18% higher. By presenting this simulation to the student, we provide concrete, visual evidence of the value of the proposed action.

## 6.4 Time-Constrained Recourse

Real-world students operate under strict time constraints. A recommendation to "study 5 hours" is useless to a student with 30 minutes before a test. Our system explicitly incorporates **Time Budget** as a constraint in the recourse optimization.

We solve for:

$$\max P(Success) \quad \text{s.t.} \quad \sum Cost(Action_i) \leq T_{available}$$

This ensures that the generated counterfactuals are not just theoretically sound but practically implementable. If the time budget is too tight to achieve the target, the system acts ethically by adjusting the goal (e.g., "Aim for a passing grade on the core sections") rather than suggesting an impossible path to a perfect score.

# 7. The Socratic Interface: Conversational Explanations

The final layer of our architecture is the interface. Raw causal graphs and SHAP values are unintelligible to most learners. To bridge the gap between algorithmic insight and human understanding, we employ a **Socratic Agent**—a Large Language Model (LLM) fine-tuned to interpret causal data and engage the student in a constructive dialogue.

## 7.1 From Data to Dialogue: Prompt Engineering

We do not simply feed the LLM the raw data. We use a structured prompt engineering strategy that serializes the SCM and counterfactual findings into a narrative context.

- **Context:** "The student failed the Unit 4 Quiz. TD-BKT estimates mastery of 'Sub-skill B' is 0.3 (Low)."
- **Causal Insight:** "SHAP analysis shows 'Lack of Practice on B' is the primary negative driver (-0.45)."
- **Counterfactual:** "Simulation shows that 15 minutes of practice on B increases pass probability to 92%."
- **Pedagogical Goal:** "Do not state the fact. Ask a question that leads the student to realize they need to practice B."

This "Chain of Thought" prompting ensures that the LLM remains grounded in the rigorous causal data while adopting the persona of a supportive tutor.

## 7.2 The Socratic Loop

Instead of a notification saying "Review Sub-skill B," the agent initiates a conversation:

- *Agent:* "I noticed the Unit 4 Quiz was tough. Looking at your practice history, we skipped the exercises for 'Sub-skill B' last Tuesday. Do you feel that gap made the quiz questions harder?"
- *Student:* "Yeah, I didn't really understand that part so I skipped it."
- *Agent:* "That makes sense. It's a tricky concept. My data suggests that just a short, 15-minute review of B could boost your confidence for the re-take significantly. Shall we look at one example together?"

This approach aligns with **Constructivist learning theory**, which posits that learners construct knowledge more effectively when they are active participants in the diagnosis of their own misconceptions. By validating the student's feeling ("Yeah, I didn't understand") with data ("We skipped exercises"), the agent builds trust.

## 7.3 Building Trust through Transparency

Trust is the currency of educational AI. If a student suspects the system is assigning "busy work," they will disengage. The "Glass Box" nature of our system mitigates this. By providing the specific rationale ("...boost your confidence significantly" backed by the 92% probability simulation), the system respects the student's time and intelligence. The Socratic dialogue serves as a negotiation interface where the user can challenge the model's assumptions (e.g., "I did study B, but I did it offline"). The agent can then update the SCM (modifying the exogenous noise variable for that study event) and recalculate the recommendations, ensuring the system remains responsive to the student's reality.

# 8. Conclusion and Future Directions

The integration of Time-Dependent Bayesian Knowledge Tracing, Multi-Armed Bandits, and Structural Causal Models represents a significant leap forward for adaptive learning systems. We move from systems that merely *predict* to systems that *explain* and *empower*.

By modeling the physics of memory through TD-BKT, we enable temporal counterfactuals that respect the reality of forgetting curves. By wrapping the optimization engine in a causal layer, we allow the system to explain its own decisions via SHAP and offer actionable recourse via counterfactual simulation. Finally, by delivering these insights through a Socratic interface, we transform algorithmic output into human-centric pedagogy.

Future work must address the challenges of "Open World" causal discovery—learning new causal links on the fly as students invent novel strategies—and the ethical dimensions of recourse. We must ensure that the "effort" required for success is distributed equitably and that the system does not inadvertently penalize students with different learning constraints. Ultimately, the goal is to create an AI that acts not as a commander, but as a wise and transparent guide, illuminating the path to mastery with the light of causal understanding.