# Behavioral Analytics and Psychometric Architectures in Adaptive Learning Systems: A Synthesis of Evidence-Centered Design, Implicit Feedback, and Developmental Measurement

## Executive Summary

The transition from static, episodic testing to dynamic, continuous assessment represents a paradigm shift in educational technology, driven by the convergence of cognitive science, machine learning, and psychometrics. Traditional assessment models, often reliant on explicit selected-response items, fail to capture the complex, temporal dynamics of learning processes, particularly in open-ended digital environments. This report provides a comprehensive technical analysis of **stealth assessment**, a methodology that embeds psychometric evaluation directly within the learning environment—such as a serious game or adaptive platform—to infer latent competencies from observable behaviors without disrupting the flow of engagement.

The architectural foundation of stealth assessment is **Evidence-Centered Design (ECD)**, a framework that establishes a rigorous evidentiary argument connecting low-level log data to high-level competency constructs. While ECD provides the inferential logic, the raw material for these inferences increasingly comes from **implicit feedback**—unsolicited behavioral signals such as click-streams, dwell times, and eye-tracking data. Interpreting these noisy signals requires sophisticated feature engineering and normalization techniques to distinguish genuine cognitive engagement from superficial interaction. Furthermore, the validity of these systems relies on **psychometric calibration** that is sensitive to the developmental stage of the learner. Standard Item Response Theory (IRT) models, particularly those accounting for guessing (3PL), often exhibit anomalous behavior when applied to pediatric populations due to systematic rather than random error patterns in children.

This report synthesizes findings from seminal projects like *Physics Playground*, *ENGAGE*, and *Mission HydroSci* to construct a unified framework for adaptive systems. It argues that the future of educational measurement lies in the move from "knowledge engineering" (expert-defined rules) to "representation learning" (deep learning-based feature extraction), and from binary learner categorization to continuous, probabilistic profiling. The analysis

highlights that while the "meshing hypothesis" of learning styles lacks empirical support, behavioral assessment of preferences remains critical for optimizing multi-modal content delivery and maintaining the Zone of Proximal Development (ZPD).

---

# 1. Stealth Assessment Implementation via Evidence-Centered Design

Stealth assessment is defined as the quiet, unobtrusive accumulation of evidence regarding a learner's competencies during the course of learning or play. Unlike traditional assessments that stop the action to test the user, stealth assessment continuously updates a learner model in real-time. The efficacy of this approach hinges on the rigorous application of Evidence-Centered Design (ECD), which prevents the "data rich, information poor" paradox by forcing designers to explicitly model the relationship between what a user *does* and what a user *knows*.

## 1.1 The Theoretical Architecture of ECD in Digital Environments

Evidence-Centered Design serves as the structural backbone for stealth assessment. It is not merely a design philosophy but a technical specification for building valid assessment engines. ECD comprises four distinct, interrelated models: the Competency Model, the Evidence Model, the Task Model, and the Assembly Model. Each model performs a specific function in the transformation of ephemeral gameplay actions into stable estimates of ability.[1]

### 1.1.1 The Competency Model (CM): Defining Latent Constructs

The Competency Model defines the "what" of the assessment—the set of Knowledge, Skills, and Abilities (KSAs) the system intends to measure. In sophisticated adaptive systems, these are not treated as binary states (Mastered/Not Mastered) but as continuous latent variables or multidimensional probability distributions.

Hierarchical and Networked Structures:
In the Physics Playground (formerly Newton's Playground) environment, the CM is structured as a Bayesian network. High-level nodes represent broad constructs such as "Newtonian Physics Understanding" or "Conscientiousness." These parent nodes branch into granular child nodes representing specific facets, such as "Newton's First Law," "Conservation of Momentum," "Properties of Torque," and "Energy Transfer".[2] This hierarchical structure allows for diagnostic granularity; a student may demonstrate high competency in "Force" but low competency in "Torque," a distinction lost in unidimensional scoring models.
Probabilistic Dependencies:
Crucially, the CM explicitly models the probabilistic dependencies between these variables. Mastery of "Force and Motion" may be modeled as a prerequisite probabilistic parent to "Energy Transfer." Mathematically, this means that an increase in the posterior probability of

the "Force" node propagates through the network, updating the prior probabilities of dependent nodes.3 This propagation allows the system to make inferences about unobserved skills based on performance in related areas, maximizing the efficiency of the assessment.

Granularity in Computational Thinking:
Effective CMs require fine-grained definitions tailored to the domain. In the ENGAGE platform, designed to teach computational thinking, the CM deconstructs broad problem-solving skills into specific algorithmic thinking components, such as "Loop Comprehension," "Conditional Logic," and "Debugging Strategy." This decomposition allows the system to track the learner's evolving mental model of programming concepts dynamically, rather than just their ability to write correct syntax.5

## 1.1.2 The Evidence Model (EM): Bridging Behavior and Construct

The Evidence Model acts as the inferential bridge between the unobservable internal states defined in the CM and the observable behaviors elicited by the system. It consists of two sub-components: **Evidence Rules** and **Statistical Models**.[2]

Evidence Rules (Feature Extraction):
Evidence rules are automated scoring rubrics that parse raw, low-level log data into "observable variables." In a physics game, raw telemetry includes data points like mouse coordinates (x, y), object creation timestamps, and simulation velocity vectors. Evidence rules translate this noise into semantic events.

- **Agent Identification:** In *Physics Playground*, a key challenge was identifying what machine a student intended to build from free-form drawings. The system utilizes an "Agent Identification System"—a rule-based classifier that analyzes the geometric properties and physics interactions of drawn objects. For instance, if an object is static and creates an inclined plane, the rule classifies it as a "ramp." If it rotates around a pivot point, it is classified as a "lever" or "pendulum".[8] This automated classification allows the system to treat unstructured creative play as structured response data.
- **Efficiency and Quality Metrics:** The rules also quantify the *quality* of the solution. The system distinguishes between "Gold Trophy" solutions (efficient, elegant solutions using fewer than three objects) and "Silver Trophy" solutions (brute-force solutions using three or more objects). This distinction is critical for the statistical model; a gold trophy provides stronger evidence of mastery and results in a larger positive update to the competency estimate than a silver trophy.[8]

Statistical Models (Scoring Engines):
The statistical model quantifies the relationship between the observable variable (e.g., drawing a lever) and the latent competency (e.g., understanding torque). Bayesian Knowledge Tracing (BKT) and Dynamic Bayesian Networks (DBN) are the standard implementation vehicles.

- **Bayesian Inference:** When a student successfully creates a functional pendulum to solve a level, this evidence is fed into the Bayes net. The "Pendulum" observable node is set to "True," which then propagates probability updates to the parent "Physics

Understanding" node based on Conditional Probability Tables (CPTs).

- **CPT Elicitation:** These CPTs define the sensitivity and specificity of the task. For example, how likely is a student with *low* physics understanding to accidentally create a working pendulum (guessing)? How likely is a student with *high* understanding to fail (slipping)? Research indicates that defining these CPT values is a recursive process involving expert elicitation and refinement using pilot data.[4]

### 1.1.3 The Task and Assembly Models

The **Task Model** specifies the features of the environment that elicit the evidence. In game-based learning, this involves designing levels that structurally require specific competencies to solve. For example, a "Pendulum Level" in *Physics Playground* is designed such that the ball cannot reach the target without the kinetic energy transfer provided by a swinging mass.[9] This ensures that success in the task is a valid proxy for the targeted skill.

The **Assembly Model** controls the adaptive sequencing of these tasks. It determines which task to present next based on the current state of the Competency Model. The goal is often to maximize information gain (reducing uncertainty about the learner's state) or to target the learner's Zone of Proximal Development (ZPD)—selecting tasks where the predicted probability of success is between 35% and 70%.[1]

## 1.2 Case Study Analysis: Physics Playground

*Physics Playground* (PP) serves as the archetype for ECD implementation in 2D simulation environments. The system utilizes a sophisticated logging mechanism that captures events at a millisecond resolution to infer both cognitive physics knowledge and non-cognitive traits like persistence and creativity.[13]

### 1.2.1 Operationalizing Persistence and Creativity

Stealth assessment in PP extends beyond domain knowledge to measure "21st-century skills."

- **Persistence Modeling:** Persistence is operationalized through behavioral markers such as the time spent on unsolved levels, the number of restarts, and the revisit rate to difficult problems. Crucially, the system must distinguish between "productive persistence" and "gaming the system" (e.g., rapid-fire guessing). The evidence rules analyze the *timing* between actions; if restarts occur too rapidly to allow for cognitive reflection, they are flagged as "gaming" and treated as negative evidence or noise.[14]
- **Creativity Metrics:** Creativity is assessed by analyzing the diversity and rarity of solutions. The Competency Model for creativity includes sub-facets like "Fluency" (number of valid ideas), "Flexibility" (variety of agent types used), and "Originality" (statistical rarity of the solution compared to the global player population). Vector-based logging tracks the trajectory of the ball; unique trajectories that solve the level are scored as highly original.[2]

### 1.2.2 The "Motif" Architecture

To scale the Bayesian network across dozens of levels without manually building a massive, unique network for the entire game, the PP team developed a "motif" architecture. A motif is a modular Bayes net fragment representing a generic task structure. When a student enters a specific level, the system instantiates the relevant motif (e.g., the "Lever Motif") and connects it to the student's persistent Competency Model. This allows for modular, scalable assessment design where new levels can be added by simply mapping them to existing motifs.[17]

## 1.3 Case Study Analysis: ENGAGE and DeepStealth

While *Physics Playground* relies on expert-defined, interpretable Bayesian networks, the *ENGAGE* platform and the *DeepStealth* framework represent a shift toward data-driven, representation learning approaches for competency estimation.

### 1.3.1 From Feature Engineering to Representation Learning

Traditional ECD relies on manual feature engineering (e.g., explicitly defining the geometry of a "ramp"). The *DeepStealth* framework, however, demonstrates that deep neural networks, specifically Long Short-Term Memory (LSTM) networks, can ingest raw game trace logs and outperform traditional Bayesian networks and Naïve Bayes classifiers in predicting student performance.[6]

- **Sequential Data Modeling:** LSTMs are particularly suited for educational data because learning is inherently temporal. The *sequence* of actions—not just the aggregate count—carries semantic weight. A student who deletes a code block and then immediately restores it is exhibiting a different cognitive process (hesitation/reflection) than one who deletes it and moves on. *DeepStealth* processes sequences of game state changes to capture these latent strategies that manual feature engineering might miss.[18]
- **Automated Feature Extraction:** In *ENGAGE*, which teaches computational thinking through a narrative-based game, the system analyzes block-based programming interactions. Rather than defining rigid rules for "loop comprehension," the deep learning model learns patterns of block manipulation (e.g., deletion rates, execution frequency, nesting depth) that correlate with post-test mastery. This approach reduces the "knowledge engineering bottleneck" associated with defining complex CPTs in Bayesian networks.[12]

### 1.3.2 Generative Zero-Shot Learning (ZSL)

A persistent challenge in data-driven assessment is the "cold start" problem—the system requires massive amounts of gameplay data to train the model for a new level. Recent advancements in *ENGAGE* and *Geniventure* (a genetics game) have introduced Zero-Shot Learning (ZSL) using Conditional Generative Adversarial Networks (CGANs).

- **Mechanism:** The CGAN is trained to generate synthetic gameplay traces that mimic human behavior for new, unseen levels based on the game's mechanics and the task

model.

- **Implication:** This allows the competency model to generalize to new content immediately. The assessment engine can predict how a competent student *would* behave in a new level and use that synthetic baseline to score real students, significantly accelerating the deployment of adaptive content.[21]

# 1.4 Algorithmic Challenges: Gaming the System and State Spaces

### 1.4.1 Detecting "Gaming the System"

"Gaming the system"—behaviors aimed at completing the task without learning (e.g., exploiting hints, systematic guessing)—is a threat to validity. In *Physics Playground*, the evidence model analyzes the inter-action latency. If a student creates and destroys objects at a frequency that exceeds human reaction time for cognitive processing, the system flags this as "off-track" behavior. In simpler systems, this might generate a "Silver Trophy" at best, but in advanced models, it triggers a "metadata" node in the Bayes net that lowers the confidence (increases variance) of the competency estimate rather than treating it as valid failure evidence.[22]

### 1.4.2 Reducing State-Space Dimensionality

In open-ended environments like *Refraction* (teaching fractions) or *DragonBox* (algebra), the state space (all possible board configurations) is theoretically infinite. Popović et al. utilized "playtraces" and state-space clustering to map common solution paths.

- **Cluster-Based Assessment:** Instead of tracking every unique coordinate, the system clusters player trajectories into "strategies." For example, in *Refraction*, players might be clustered into "Optimal Splitters" (who use the most efficient fraction dividers) vs. "Recursive Error Makers."
- **Constraint-Based Modeling:** This reduces the dimensionality of the assessment problem. The system only needs to identify which cluster a student's current trajectory belongs to in order to predict their final outcome and intervene if they are on a known "failure path".[23]

# 1.5 Table: Comparison of Assessment Frameworks

| Feature | Physics Playground (ECD/Bayes) | ENGAGE / DeepStealth (Deep Learning) |
|---|---|---|
| **Core Architecture** | Dynamic Bayesian Networks (DBN) | Long Short-Term Memory (LSTM) / CNN |
| **Feature Extraction** | Manual (Rule-based Agent | Automated (Representation |

|  | ID) | Learning) |
|---|---|---|
| **Interpretability** | High (Explicit nodes & CPTs) | Low (Black-box hidden layers) |
| **Data Requirement** | Low (Expert priors can bootstrap) | High (Requires massive training sets) |
| **Scalability** | Linear (New motifs for new levels) | High (Zero-Shot Learning for new levels) |
| **constructs** | Physics, Persistence, Creativity | Computational Thinking, Problem Solving |

---

# 2. Implicit Feedback Interpretation in Educational Technology

As adaptive systems evolve, they increasingly rely on **implicit feedback**—unsolicited signals generated by the user's natural interaction with the system—to refine learner profiles. Unlike explicit feedback (e.g., quiz scores, self-reported surveys), which is sparse and subjective, implicit feedback (e.g., dwell time, click patterns, gaze) is abundant and objective. However, interpreting this data requires sophisticated normalization and fusion techniques to filter out noise.

## 2.1 Click-Stream Analysis and Preference Inference

Click-stream data provides the foundational layer of implicit feedback. However, raw click counts are often misleading indicators of learning or preference due to "clickbait" effects, navigational errors, or random exploration.

### 2.1.1 Relative Preference Valuation

A cornerstone finding in implicit feedback research, established by Joachims et al., is that **relative preferences are significantly more reliable than absolute judgments**.

- **The "Click > Skip Above" Heuristic:** In a list of learning resources (or search results), if a user skips item A to click item B, it can be inferred with high confidence that $Preference(B) > Preference(A)$, even if the absolute relevance of B is unknown. This eliminates the need to determine a universal "relevance score" for every action.[1]
- **Educational Application:** If a learner consistently skips text-heavy modules to engage with video simulations, the system infers a relative modality preference ($Video > Text$).

This signal is stronger than simply observing that the user watched a video, as the *rejection* of the alternative provides critical context.[1]

### 2.1.2 Click-Through Rate (CTR) and Validity

While CTR is a standard metric in web analytics, in education, a "click" does not equal "learning." The interpretation must differentiate between "navigational clicks" (searching for content) and "engagement clicks" (consuming content).

- **Valid Read Detection:** Systems must filter for "valid reads"—interactions where the dwell time exceeds a cognitive threshold required for information processing. Xie et al. demonstrated that reweighting clicks based on dwell time significantly improves recommendation accuracy by filtering out "clickbait" or accidental clicks.[26]

## 2.2 Dwell Time: Normalization and Semantic Meaning

Dwell time is a high-bandwidth signal, but it is highly context-dependent. A 30-second dwell on a complex diagram suggests deep engagement; 30 seconds on a simple sentence suggests confusion or disengagement (absence).

### 2.2.1 Normalization Techniques

To use dwell time as a valid proxy for interest or difficulty, it must be normalized against the content's inherent characteristics.

- Text Length Normalization: A standard approach is to normalize dwell time ($DT$) by the number of words ($N$) or the complexity of the resource. The formulaic approach often involves calculating a "reading rate" ($R = N / DT$) and comparing it to average user baselines ($R_{avg}$). A dwell time is considered significant only if:

  $$DT > \alpha \cdot \left( \frac{N}{R_{avg}} \right)$$

  where $\alpha$ is a content-specific coefficient accounting for difficulty.28
- **Log-Normal Distribution:** Dwell time data typically follows a log-normal distribution (right-skewed). Statistical models must apply a logarithmic transformation ($\ln(DT)$) before using the data in linear regressions or clustering algorithms. This prevents outliers (e.g., a user leaving the browser open while taking a break) from skewing the learner model.[30]
- **Relative Dwell Time (RDT):** In visual tasks, RDT is calculated as the time spent on a specific Area of Interest (AOI) divided by the total time on the task. High RDT on relevant AOIs correlates with expertise, while high RDT on irrelevant AOIs (distractors) is a strong predictor of misconceptions or low competency.[31]

### 2.2.2 Contextual Semantics: Confusion vs. Engagement

The semantic meaning of dwell time flips depending on the task type.

- **Consumption Tasks (Reading/Watching):** In reading tasks, high dwell time generally correlates with interest and positive preference.[33] However, excessively high dwell time combined with regression (eye-movement backtracking) can indicate reading difficulty, dyslexia, or low fluency.[34]
- **Problem-Solving Tasks:** In environments like *DragonBox*, high dwell time on a specific step often signals "confusion" or "cognitive disequilibrium." The interpretation depends on the *subsequent action*:
  - **High Dwell $\rightarrow$ Correct Action:** Indicates productive struggle and successful engagement with the ZPD.
  - **High Dwell $\rightarrow$ Incorrect Action/Exit:** Indicates frustration and failure. This pattern triggers scaffolding or hints in adaptive systems.[36]

## 2.3 Multi-Modal Signal Fusion

The most robust learner models are constructed by fusing multiple data streams. **Multi-Modal Learning Analytics (MMLA)** integrates log data with physiological signals (eye-tracking, EEG, facial expression) to triangulate cognitive states.[37]

### 2.3.1 Fusion Architectures

- **Early Fusion (Feature Level):** Features from different modalities (e.g., fixation duration from eye-tracking and click counts from logs) are concatenated into a single vector *before* being fed into a classifier. This is effective when the modalities are highly correlated and synchronous.[38]
- **Late Fusion (Decision Level):** Separate models are trained for each modality (e.g., a CNN for facial expression and an LSTM for log sequences). Their independent predictions (e.g., "Frustrated" vs. "Engaged") are then aggregated using weighted voting or a meta-classifier. Late fusion is generally superior for educational data because it handles missing data (e.g., camera obstruction, sensor noise) more gracefully and allows for modality-specific timescales.[38]
- **Hybrid/Deep Fusion:** Recent approaches use deep neural networks where modalities interact at intermediate layers. For instance, an attention mechanism might weigh the "eye-tracking" vector more heavily during reading tasks and the "click-stream" vector more heavily during interactive tasks, dynamically adjusting the influence of each modality based on the context.[40]

### 2.3.2 Eye-Tracking as a Learning Style Indicator

Eye-tracking provides the most direct physiological evidence of learning style preferences, validating the behavioral indicators derived from logs.

- **Visual Learners:** Exhibit distinctive scan paths—they fixate on diagrams/images first (shorter time-to-first-fixation on visuals) and have higher RDT on graphical AOIs.[1]
- **Textual Learners:** Show systematic linear scanning of text blocks and lower RDT on

decorative or supplementary graphics.[1]

- **Analysis of Moving AOIs:** In dynamic games, tracking attention is computationally expensive because objects move. New methods using Vision Transformers (ViT) have been developed to automatically map gaze to moving objects (dynamic AOIs) without frame-by-frame manual annotation, achieving >99% accuracy in hold-out tests. This technology allows adaptive games to know exactly which moving game asset a student is tracking, providing granular data on visual attention strategies.[42]

## 2.4 Continuous Profiling vs. Binary Classification

The integration of implicit feedback supports the move from binary learning style labels (e.g., "Visual Learner" vs. "Auditory Learner") to **continuous probability distributions**.

Percentage-Based Profiles:
Instead of categorizing a user as "Visual," the system maintains a preference vector, such as:

$$Profile_{User} =$$

For example: $[0.45, 0.20, 0.25, 0.10]$. This vector represents the probability that the user will engage most effectively with content of that modality.[1]
Bayesian Updating with Exponential Decay:
As new implicit evidence $E$ is observed (e.g., user skips text, plays video), the profile is updated using Bayes' rule:

$$P(Style|E) \propto P(E|Style) \cdot P(Style)$$

To account for evolving preferences and the Expertise Reversal Effect (where preferences change as learners become experts, often moving from visual scaffolding to textual/symbolic efficiency), recent observations are weighted more heavily using an exponential decay function on historical data. This ensures the profile remains current and responsive to the learner's development.[1]

---

# 3. Psychometric Properties of Adaptive Assessment Across Developmental Stages

The validity of stealth assessment and adaptive algorithms relies on psychometric models, primarily **Item Response Theory (IRT)**. However, a critical failure mode in EdTech is the uncritical application of adult-normed psychometric assumptions to children. The cognitive and behavioral architecture of children necessitates specific calibrations of IRT parameters,

particularly regarding guessing and slipping behaviors.

## 3.1 Item Response Theory (IRT) Fundamentals in Adaptive Contexts

IRT models the probability $P(\theta)$ of a correct response as a function of the learner's latent ability ($\theta$) and item parameters. The standard models used in adaptive testing are:

- **1PL (Rasch Model):** Models only item difficulty ($b$).
- **2PL Model:** Adds item discrimination ($a$).
- 3PL Model: Adds a pseudo-guessing parameter ($c$).

$$P(\theta) = c + (1 - c) \frac{1}{1 + e^{-a(\theta - b)}}$$

## 3.2 The Guessing Parameter (c-parameter) Anomaly in Children

A pervasive finding in the research is that the **3PL model's guessing parameter operates differently for young children compared to adults**, often rendering standard 3PL models invalid for K-12 assessment.

### 3.2.1 Random vs. Systematic Error

The $c$-parameter assumes that low-ability examinees guess randomly (e.g., a 25% chance on a 4-option multiple choice). However, young children often do not guess randomly. They exhibit **systematic error patterns** driven by developmental misconceptions, attractors (e.g., picking the prettiest picture), or perseveration (picking the same position repeatedly).[43]

- **Empirical Evidence:** Studies on instruments like the *Children's Behavior Questionnaire* indicate that for young populations (e.g., grades K-2), the inclusion of a guessing parameter often leads to convergence failures or instability in parameter estimation because the "guessing" is not stochastic.[45] The behaviors are deterministic based on non-construct factors.
- **Recommendation:** For assessments targeting early childhood, **1PL (Rasch) or 2PL models are often psychometrically superior to 3PL models**. The 1PL model's insensitivity to guessing requires the interface design to minimize guessing opportunities (e.g., using constructed response interactions like dragging and dropping, rather than multiple choice), but it provides more stable theta estimation for developmental cohorts.[45]

## 3.3 The "Slipping" Parameter (d-parameter) and the 4PL Model

While the guessing parameter ($c$) handles low-ability students getting lucky, the **slipping parameter ($d$)** handles high-ability students making careless errors. The 4PL model adds an upper asymptote ($d < 1$).

### 3.3.1 Relevance for Pediatric Populations

The 4PL model is particularly relevant for children. Research on Human Figure Drawings (HFD) and other performance tasks shows that high-ability children are more prone to "slipping" than high-ability adults. This is often due to developmental constraints in fine motor skills, attention span, or impulsivity, rather than a lack of cognitive mastery.[48]

- **Implication:** A high-ability child might know the answer but fail to click the correct object due to motor inaccuracy. A standard 2PL or 3PL model would penalize this heavily, lowering the $\theta$ estimate. A 4PL model accounts for this non-zero probability of failure at high ability levels, providing a "fairer" assessment of cognitive competence distinct from motor performance.[48]

## 3.4 Differential Item Functioning (DIF) Across Ages

Developmental stages introduce non-construct variance that can bias assessment. An item might be difficult for a 7-year-old not because of the math content (the construct), but because of the reading load or fine motor requirements.

- **DIF Analysis:** IRT analysis must rigorously test for Differential Item Functioning across age groups. Items that show high DIF—meaning they function differently for children of the same ability level but different ages—must be flagged. This ensures that the assessment measures the intended construct ($\theta$) rather than developmental maturity or reading fluency.[50]

## 3.5 Validity of Continuous Proficiency Estimation

Adaptive systems perform "continuous assessment," updating $\theta$ in real-time. The validity of this approach depends on the stability of the construct over time.

### 3.5.1 Temporal Granularity and Latent Trait Stability

Standard IRT assumes $\theta$ is constant during the test. In learning games, $\theta$ changes *by design* because learning is occurring. Therefore, models must evolve from static IRT to **Dynamic Bayesian Networks (DBN)** or **T-SKIRT** (Temporal Structured Knowledge IRT). These models explicitly model the transition of $\theta$ from time $t$ to $t+1$, treating ability as a dynamic state variable rather than a fixed trait.[51]

### 3.5.2 Time-Series Validation Strategies

Validating these continuous models requires specialized techniques. Standard k-fold cross-validation is invalid because it shuffles future data into the training set, causing data leakage.

- **Rolling Origin Validation:** Validation must use **time-series cross-validation** (also known as "rolling origin" or "walk-forward" validation). The model is trained on data from $t_{0}...t$ and tested on $t+1$. This mirrors the real-world operation of the adaptive engine and prevents the model from "peeking" at future learning outcomes to predict

current performance.[52]

### 3.5.3 Predictive Validity

The "gold standard" for these systems is predictive validity against external standardized measures.

- **Physics Playground:** Competency estimates derived from the game's Bayesian network showed significant correlation ($r \approx 0.40 - 0.60$) with external paper-and-pencil physics pre/post-tests. This establishes concurrent validity—the game measures the same construct as the test.[8]
- **ENGAGE:** Behavioral sequences clustered by LSTM models successfully predicted post-test performance with higher accuracy than static measures. This demonstrates that the *process* of solving (strategy) is as predictive, if not more so, than the correctness of the solution itself.[18]

---

# 4. Technical Architecture and Implementation Recommendations

The convergence of ECD, implicit feedback analysis, and developmental psychometrics points toward a unified architecture for next-generation educational tools. This section outlines the technical requirements for implementing such a system.

## 4.1 System Design and Data Pipelines

Implementing stealth assessment requires a microservices architecture capable of real-time processing.

- **Latency Requirements:** Real-time adaptation (e.g., generating a hint or adjusting difficulty) requires decision latency under 200ms. This necessitates edge computing or highly optimized inference engines for the Bayesian/LSTM models.[54]
- **Data Pipeline (IDEFA Framework):** The **Integrated Design of Event-stream Features for Analysis (IDEFA)** framework suggests a rigorous pipeline:
  1. **Raw Log Ingestion:** Capture low-level events (clicks, mouse moves) at millisecond resolution.
  2. **Feature Engineering:** Aggregate raw logs into "Base Features" (e.g., dwell time, object count).
  3. **Intermediate Feature Generation:** Apply ECD rules (e.g., classify object as "ramp").
  4. **Competency Estimation:** Feed features into the DBN or LSTM to update $\theta$.
  5. **Adaptive Actuation:** The Assembly Model selects the next task based on the new $\theta$.[55]

## 4.2 Table: Architectural Requirements by Component

| Component | Standard Approach (Legacy) | Adaptive / Stealth Approach (Recommended) | Technical Implementation Strategy |
|---|---|---|---|
| **Assessment Model** | Explicit Testing (Quizzes) | Stealth Assessment (Behavioral) | **ECD Framework**: Use Bayesian Networks for interpretability (Physics Playground) or LSTMs for raw log ingestion (ENGAGE). |
| **Learner Profiling** | Binary Labels (e.g., "Visual Learner") | Continuous / Multimodal Profiles | **Vector-based profiles**: $[P_{vis}, P_{aud}, P_{kin}]$. Use Bayesian updating with exponential decay for recent behaviors. |
| **Feedback Signal** | Correct / Incorrect Answer | Implicit Signals (Dwell, Click patterns) | **Normalization**: Dwell time normalized by reading rate ($N/R_{avg}$). **Relative Preference**: Click > Skip logic. |
| **Fusion Strategy** | Unimodal (Score only) | Multi-modal (Logs + Eye/Physio) | **Late Fusion**: Ensemble classifiers to handle missing sensor data robustly. |
| **Psychometrics** | 3PL IRT (One size | Developmental IRT | **Model Selection**: Use 1PL/2PL for |

|  | fits all) | (1PL/2PL/4PL) | young children to avoid guessing instability; use 4PL to account for slipping/motor errors. |
|---|---|---|---|
| **Validation** | K-Fold Cross-Validation | Time-Series Validation | **Walk-Forward**: Train on $t_{0}..t$, predict $t+1$. |

---

# 5. Strategic Synthesis and Future Directions

The integration of stealth assessment, implicit feedback, and developmental psychometrics represents the maturation of educational data science. We are moving from systems that "test" to systems that "sense."

**Key Strategic Insights:**

1. **Behavior Over Self-Report:** The data decisively supports utilizing behavioral extraction over self-reported surveys for learning styles. Users often lack the metacognitive awareness to accurately report their preferences, but their click-stream and dwell patterns reveal their true "functioning" style. The "meshing hypothesis" may be invalid, but offering *variety* based on preference is essential for engagement.[1]
2. **The "Slipping" Parameter is Critical for Children:** While the "guessing" ($c$) parameter is unstable in young children, the "slipping" ($d$) parameter is essential. Children often know the answer but fail due to motor/attention issues. Adaptive systems for K-8 must account for this to avoid systematically under-estimating competency.
3. **Process is the Product:** In systems like *ENGAGE* and *DragonBox*, the *sequence* of actions (the strategy) is a more potent predictor of mastery than the final answer. Feature engineering must focus on sequential patterns (e.g., "how did they arrive at the solution?") rather than static states.
4. **Relative Preference Reliability:** The absolute time spent on a resource is noisy; the ratio of time spent on Resource A vs. Resource B within the same session is a high-fidelity signal. Algorithms should optimize for ranking relative preferences rather than predicting absolute ratings.

Conclusion:
Building effective adaptive educational tools requires a rigorous synthesis of cognitive science and data engineering. The implementation of stealth assessment using ECD provides the theoretical validity needed to trust the system's inferences. However, these inferences are

only as good as the raw data interpretation—necessitating advanced normalization of implicit feedback and multi-modal fusion. Finally, the "human" element—the developmental stage of the learner—dictates the statistical boundaries of the system. We cannot simply shrink adult psychometric models to fit children; we must calibrate the mathematical architecture to reflect the unique cognitive and behavioral realities of the developing mind. The future of adaptive learning lies in systems that facilitate a continuous, invisible, and developmentally calibrated dialogue with the learner.

## Works cited

1. Building Adaptive Educational Tools_ A Comprehensive Framework for Behavioral Learning Assessment.pdf
2. Stealth assessment: A theoretically grounded and psychometrically sound method to assess, support, and investigate learning in technology-rich environments, accessed January 5, 2026, https://myweb.fsu.edu/vshute/pdf/ETRD2023.pdf
3. Stealth assessment in computer-based games to support learning - Florida State University, accessed January 5, 2026, https://myweb.fsu.edu/vshute/pdf/shute%20pres_h.pdf
4. An Elicitation Tool for Conditional Probability Tables (CPT) for Physics Playground, accessed January 5, 2026, https://www.researchgate.net/publication/321316020_An_Elicitation_Tool_for_Conditional_Probability_Tables_CPT_for_Physics_Playground
5. Inducing Stealth Assessors from Game Interaction Data, accessed January 5, 2026, https://intellimedia.ncsu.edu/wp-content/uploads/sites/42/min-aied-2017.pdf
6. DEEPSTEALTH: Game-Based Learning Stealth Assessment with Deep Neural NetworCs - IEEE Xplore, accessed January 5, 2026, https://ieeexplore.ieee.org/ielaam/4620076/9120392/8735739-aam.pdf
7. (PDF) Stealth Assessment - ResearchGate, accessed January 5, 2026, https://www.researchgate.net/publication/363019884_Stealth_Assessment
8. Assessment and Learning of Qualitative Physics in Newton's Playground - Florida State University, accessed January 5, 2026, https://myweb.fsu.edu/vshute/pdf/JER.pdf
9. How Task Features Impact Evidence From Assessments Embedded in Simulations and Games - PMC - NIH, accessed January 5, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC6176773/
10. Final version of Bayesian net for PP | Download Scientific Diagram - ResearchGate, accessed January 5, 2026, https://www.researchgate.net/figure/Final-version-of-Bayesian-net-for-PP_fig2_321316020
11. Newton's Playground: How to use evidence centered design (ECD) to develop game-based assessment - ETC Press, accessed January 5, 2026, https://press.etc.cmu.edu/file/download/1594/8fa419b2-2531-4364-96b3-0bb8254ad719

12. Stealth Assessment - Valerie Shute, Xi Lu and Seyedahmad Rahimi - ERIC, accessed January 5, 2026, https://files.eric.ed.gov/fulltext/ED612156.pdf
13. A snapshot of a Physics Playground level log file - ResearchGate, accessed January 5, 2026, https://www.researchgate.net/figure/A-snapshot-of-a-Physics-Playground-level-log-file_fig10_282892197
14. Construct and Predictive Validity of an Assessment Game to Measure Honesty–Humility, accessed January 5, 2026, https://www.researchgate.net/publication/348515581_Construct_and_Predictive_Validity_of_an_Assessment_Game_to_Measure_Honesty-Humility
15. Iterative Feature Engineering Through Text Replays of Model Errors - Penn Center for Learning Analytics, accessed January 5, 2026, https://learninganalytics.upenn.edu/ryanbaker/paper%20141b.pdf
16. (PDF) Stealth Assessment and Digital Learning Game Design - ResearchGate, accessed January 5, 2026, https://www.researchgate.net/publication/375223668_Stealth_Assessment_and_Digital_Learning_Game_Design
17. Applying Evidence-Centered Design for the Development of Game-Based Assessments in Physics Playground - Florida State University, accessed January 5, 2026, https://myweb.fsu.edu/vshute/pdf/IJT.pdf
18. DeepStealth: Game-Based Learning Stealth Assessment With Deep Neural Networks, accessed January 5, 2026, https://ieeexplore.ieee.org/document/8735739/
19. Improving Stealth Assessment in Game-based Learning with LSTM-based Analytics - ERIC, accessed January 5, 2026, https://files.eric.ed.gov/fulltext/ED593099.pdf
20. Game-Based Learning Prediction Model Construction | Journal of Learning Analytics, accessed January 5, 2026, https://learning-analytics.info/index.php/JLA/article/view/8105
21. Enhancing Stealth Assessment in Game-Based Learning Environments with Generative Zero-Shot Learning - Educational Data Mining, accessed January 5, 2026, https://educationaldatamining.org/edm2022/proceedings/2022.EDM-long-papers.15/index.html
22. Debugging the Evidence Chain - CEUR-WS.org, accessed January 5, 2026, https://ceur-ws.org/Vol-1024/paper-01.pdf
23. A Data-Driven Approach for Inferring Student Proficiency from Game Activity Logs - University of Pittsburgh, accessed January 5, 2026, https://people.cs.pitt.edu/~falakmasir/docs/gain_paper2016.pdf
24. Feature-Based Projections for Effective Playtrace Analysis - University of Washington, accessed January 5, 2026, https://grail.cs.washington.edu/projects/playtracer/fdg2011/fdg2011.pdf
25. Evaluating the Accuracy of Implicit Feedback from Clicks and Query Reformulations in Web Search - Cornell: Computer Science, accessed January 5, 2026, https://www.cs.cornell.edu/~tj/publications/joachims_etal_07a.pdf

26. Reweighting Clicks with Dwell Time in Recommendation - arXiv, accessed January 5, 2026, https://arxiv.org/pdf/2209.09000
27. Beyond Explicit and Implicit: How Users Provide Feedback to Shape Personalized Recommendation Content - arXiv, accessed January 5, 2026, https://arxiv.org/html/2502.09869v1
28. Feature Selection and Model Comparison on Microsoft Learning-to-Rank Data Sets - arXiv, accessed January 5, 2026, https://arxiv.org/pdf/1803.05127
29. Lessons from the Journey: A Query Log Analysis of Within-Session Learning - Microsoft, accessed January 5, 2026, https://www.microsoft.com/en-us/research/wp-content/uploads/2016/04/wsdm14.pdf
30. Streaming, Fast and Slow: Cognitive Load-Aware Streaming for Efficient LLM Serving - arXiv, accessed January 5, 2026, https://arxiv.org/html/2504.17999v1
31. A Literature Review Comparing Experts' and Non-Experts' Visual Processing of Graphs during Problem-Solving and Learning - MDPI, accessed January 5, 2026, https://www.mdpi.com/2227-7102/13/2/216
32. Teachers' gaze over space and time in a real-world classroom - PMC - NIH, accessed January 5, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC8003357/
33. 38 Towards Individuated Reading Experiences: Different Fonts Increase Reading Speed for Different Individuals - Jeff Huang, accessed January 5, 2026, https://jeffhuang.com/papers/Readability_TOCHI22.pdf
34. Reading Specific Small Saccades Predict Individual Phonemic Awareness and Reading Speed - Frontiers, accessed January 5, 2026, https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2021.663242/full
35. Individuals with dyslexia use a different visual sampling strategy to read text - PMC, accessed January 5, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC7979812/
36. An Exploratory Analysis of Confusion Among Students Using Newton's Playground - Penn Center for Learning Analytics, accessed January 5, 2026, https://learninganalytics.upenn.edu/ryanbaker/ICCE%20NP%20v08.pdf
37. Multimodal Data Fusion in Learning Analytics: A Systematic Review - MDPI, accessed January 5, 2026, https://www.mdpi.com/1424-8220/20/23/6856
38. Effective Techniques for Multimodal Data Fusion: A Comparative Analysis - PMC, accessed January 5, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC10007548/
39. Multi-physiological signal fusion for objective emotion recognition in educational human–computer interaction - ResearchGate, accessed January 5, 2026, https://www.researchgate.net/publication/386581572_Multi-physiological_signal_fusion_for_objective_emotion_recognition_in_educational_human-computer_interaction
40. Emotion recognition and achievement prediction for foreign language learners under the background of network teaching - NIH, accessed January 5, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC9637875/
41. Consistent but modest: A meta-analysis on unimodal and multimodal affect detection accuracies from 30 studies - ResearchGate, accessed January 5, 2026,

https://www.researchgate.net/publication/266652990_Consistent_but_modest_A_meta-analysis_on_unimodal_and_multimodal_affect_detection_accuracies_from_30_studies

42. Quantifying Dwell Time With Location-based Augmented Reality: Dynamic AOI Analysis on Mobile Eye Tracking Data With Vision Transformer - NIH, accessed January 5, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC11165940/

43. Item Response Theory Analyses of the Cambridge Face Memory Test (CFMT), accessed January 5, 2026, https://www.researchgate.net/publication/271706135_Item_response_theory_analyses_of_the_Cambridge_Face_Memory_Test_CFMT

44. Development of a Computerized Adaptive Test in Human Growth and Development to Enhance Counselor Preparation Comprehensive Exami - Scholarship@Miami, accessed January 5, 2026, https://scholarship.miami.edu/view/pdfCoverPage?instCode=01UOML_INST&filePid=13450821090002976&download=true

45. Examining the Measurement Precision and Invariance of the Revised Get Ready to Read! - PMC - NIH, accessed January 5, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC4762015/

46. Early Childhood Longitudinal Study, Birth Cohort (ECLS–B) Psychometric Report for the 2-Year Data Collection - National Center for Education Statistics (NCES), accessed January 5, 2026, https://nces.ed.gov/pubs2007/2007084_C2.pdf

47. Full article: The youth social media literacy inventory: Development and validation using item response theory in the US, accessed January 5, 2026, https://www.tandfonline.com/doi/full/10.1080/17482798.2023.2230493

48. Using four-parameter item response theory to model human figure drawings - Pepsic, accessed January 5, 2026, https://pepsic.bvsalud.org/scielo.php?script=sci_arttext&pid=S1677-04712018000400008

49. Using Four-Parameter Item Response Theory to model Human Figure Drawings, accessed January 5, 2026, https://www.researchgate.net/publication/332857723_Using_Four-Parameter_Item_Response_Theory_to_model_Human_Figure_Drawings

50. Psychometric Properties and Normative Data Using Item Response Theory Approach for Three Neuropsychological Tests in Waranka Children Population - MDPI, accessed January 5, 2026, https://www.mdpi.com/2227-9032/13/4/423

51. T-SKIRT: Online Estimation of Student Proficiency in an Adaptive Learning System, accessed January 5, 2026, https://www.researchgate.net/publication/313712211_T-SKIRT_Online_Estimation_of_Student_Proficiency_in_an_Adaptive_Learning_System

52. AI for Data Analysis | Cross-Validation - Julius AI, accessed January 5, 2026, https://julius.ai/glossary/cross-validation

53. What Is Cross-Validation? A Plain English Guide with Diagrams - KDnuggets, accessed January 5, 2026, https://www.kdnuggets.com/what-is-cross-validation-a-plain-english-guide-with-diagrams

54. Challenging Cognitive Load Theory: The Role of Educational Neuroscience and Artificial Intelligence in Redefining Learning Efficacy - PubMed Central, accessed January 5, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC11852728/
55. Fueling Prediction of Player Decisions: Foundations of Feature Engineering for Optimized Behavior Modeling in Serious Games, accessed January 5, 2026, https://learninganalytics.upenn.edu/ryanbaker/owen-tknl.pdf
56. Game-Based Learning Prediction Model Construction: Toward Validated Stealth Assessment Implementation - ERIC, accessed January 5, 2026, https://files.eric.ed.gov/fulltext/EJ1465546.pdf