



Proficiency Gain Estimation from Non-Question Learning Activities – Comprehensive Research Report

Executive Summary

Overview: Developing a proficiency evaluation framework that incorporates **non-question learning activities** (e.g. reading texts, watching videos, interactive simulations) is crucial for Pearson's Nexus AI platform. Traditional assessments rely on quiz or test questions to gauge learning, but modern learning platforms gather rich **engagement data** from content interactions. This report surveys literature and industry practices to identify how such data can be used to **infer proficiency gains**, and recommends a multi-faceted framework applicable across ages and subjects.

Key Findings:

- **Learning Engagement Correlates with Outcomes:** Research consistently shows that higher engagement with learning content (time spent, completeness, re-engagement behaviors) correlates with better learning outcomes ¹ ². For example, *every additional minute on Khan Academy* was associated with gains on standardized tests ¹, and **careful video viewing behaviors** (pausing, rewinding) correlate with higher exam performance ³. While correlation isn't causation, these patterns provide evidence that *meaningful content interaction yields proficiency gains*.
- **Challenges of Non-Question Assessment:** Unlike quizzes, content interactions don't yield a clear "right or wrong" result. This makes proficiency inference an **estimation problem under uncertainty**. Cognitive science warns that passive study (e.g. reading or watching) is less predictive of retention than active practice. Thus, models must account for **diminished or variable learning gains** from content consumption. Key challenges include:
 - **Detecting Actual Engagement:** Simply logging time is insufficient – a student might spend 10 minutes on a page but learn nothing if distracted. Fine-grained metrics (scroll depth, video play/pause events, note-taking) are needed to gauge genuine engagement.
 - **Variability in Learning Efficacy:** Two students may spend the same time on a video but learn differently. Prior knowledge, learning skills, and content difficulty mediate proficiency gain.
 - **Lack of Immediate Feedback:** Without question responses, the system must infer learning indirectly, requiring **probabilistic models** or delayed validation via later assessments.
- **Methodologies for Estimating Proficiency Gains:**
 - **Heuristic Point Systems (XP Models):** Many platforms assign experience points or progress percentages for completing content. This is simple and age-agnostic: e.g. Duolingo awards XP for

lessons (regardless of whether a quiz question was answered) to encourage continued engagement. However, XP is more a **gamification metric** than a true knowledge measure – it assumes any engagement is positive. While easy to implement, this approach must be calibrated to avoid misleading high-engagement/low-mastery scenarios.

- **Time-on-Task & Completion Metrics:** A refined approach uses **time-based weighting** and content completion status. For example, MOOCs use *video attendance rate* (% of videos watched) and *utilization rate* (% of video duration watched) as predictors of course success ⁴ ⁵. A simple model might award partial mastery credit for spending sufficient time on a content item. This can be formulated as:

$\Delta \text{Proficiency}_{\text{content}} = w_c \times f(\text{engagement})$
where w_c is a weight for the content's instructional value, and $f(\text{engagement})$ is a function of the learner's engagement (e.g. fraction of content completed, normalized time spent). Studies show time and completion metrics have predictive power – e.g. *attendance and utilization rates in week 1 can highly predict which students will pass a MOOC* ⁴. Nevertheless, raw time must be adjusted for quality: **active engagement behaviors** contribute more than passive time. For instance, a study found **number of pauses and rewinds in videos was positively correlated with test scores**, whereas excessive fast-forwarding correlated negatively ³.

- **Behavioral Analytics & “Stealth Assessment”:** Advanced techniques analyze **interaction patterns** (beyond duration) to infer skill mastery. In educational games and interactive simulations, every action (moves, errors, hints) can be evidence of proficiency. Pearson's own research introduced *SPRING (Student PRoficiency INFerrer from Game data)* – a data-driven pipeline that learned to predict test outcomes from game logs ⁶ ⁷. Without any direct quiz questions, SPRING achieved a correlation ~0.55 with external assessment scores, validating that **sequence-of-actions data can predict learning outcomes** ⁶. The **Evidence-Centered Design (ECD)** framework underpins such stealth assessment: define a competency model (skills to measure) and an evidence model mapping behaviors to skill mastery ⁸ ⁹. Historically this required expert-crafted rules (e.g. Bayesian networks in Shute's stealth assessments), but modern approaches leverage machine learning to discover patterns automatically ¹⁰ ¹¹. These methods are **content-agnostic** (can be applied to math games or history simulations alike) and age-neutral, focusing on how the student learns rather than the specific subject content.

- **Knowledge Modeling (Bayesian/IRT Approaches):** In adaptive learning systems, it's common to maintain a **probability of mastery** for each skill (as in Bayesian Knowledge Tracing). Typically this updates when a student answers a question correctly or incorrectly. To incorporate content interactions, we treat them as *learning opportunities* with uncertain outcomes. For example, upon **reading a chapter or watching a tutorial**, the system can increase the estimated mastery of relevant skills by a certain amount (less than if the student had answered practice questions). A possible Bayesian update rule:

$$P(\text{mastery} | \text{content}) = P(\text{mastery before}) + (1 - P(\text{mastery before})) \times L_c$$

where L_c is the **learning gain parameter** for that content piece (akin to a “probability of learning” from the material). This formula (inspired by knowledge-tracing models) means if a student hasn't mastered a skill yet, engaging with content gives them some chance L_c to learn it ¹². The L_c can be calibrated – for example, reading an explanatory text might have a lower L_c than solving a problem. If the student was already near mastery, the gain is minimal (ceiling effect). Over time, as the system collects data, it can adjust L_c for each content type by comparing predicted vs. actual performance (i.e. if many students still fail a skill after just reading, the L_c for reading alone should be set low).

- **Machine Learning Predictive Models:** With sufficient data, one can train regression or classification models to *directly predict proficiency* (or test scores) from features of content interaction logs. For example, researchers have used **clickstream features** (e.g. number of forum posts, videos watched, assignments opened) to predict final exam scores in MOOCs. In one study, a combination of video interaction features predicted student grades with reasonable accuracy early in the course ⁴. These models often use algorithms like random forests or neural networks to capture nonlinear relationships among behaviors. A recent direction is **Deep Knowledge Tracing (DKT)**, a neural network that processes sequences of student interactions to predict future performance. DKT traditionally uses question response sequences; however, research is emerging to integrate non-question data – e.g. **DKT with cognitive load estimation** ¹³ ¹⁴. By incorporating indicators of mental effort (from timing, eye-tracking, etc.), such models aim to more holistically estimate learning state, not just whether a question was answered correctly. While powerful, ML models require large datasets and careful validation to ensure they generalize and provide interpretable insights.

Industry Case Studies:

- **Khan Academy:** Khan Academy's mastery system is primarily based on **practice exercises and quizzes**. In fact, *watching videos or reading articles does not directly count toward "mastery" in Khan's platform* ¹⁵. Videos are seen as support tools; mastery is demonstrated by doing. This conservative approach avoids over-estimating proficiency from passive learning. However, Khan Academy does track content engagement for **analytics purposes** – teachers and students can see if videos were watched, and Khan's research has shown time spent on the platform correlates with external test score gains ¹. The takeaway is that Khan Academy prioritizes **assessment-based evidence** for formal proficiency but acknowledges content engagement as an important supporting factor (especially in driving practice: videos help students eventually solve problems).
- **Duolingo:** Duolingo offers an interesting two-layer model. On the surface, it uses **XP points and streaks** to encourage daily engagement (a motivator, not a direct proficiency metric). Under the hood, Duolingo pioneered a data-informed proficiency model called **Half-Life Regression (HLR)** ¹⁶. HLR is a statistical model that treats each practice event (translating a phrase, etc.) as evidence and fits a **forgetting curve** for each word or skill. The model predicts the probability the learner can recall a word at any given time, based on practice history and spacing ¹⁶ ¹⁷. This drives the "strength bars" users see for skills. While HLR relies on question attempts (translations) as inputs, its methodology is relevant: it transforms *behavior data over time into an estimate of current proficiency*. Duolingo doesn't award proficiency credit for simply reading a lesson tip or listening to a phrase; it always pairs content with an interactive prompt. The lesson for our framework is the power of combining **learning science (forgetting curves)** with big data to model proficiency continually. Even if Nexus AI includes pure content interactions, we can adapt Duolingo's approach: treat content interaction as a less-potent practice event and update a forgetting curve or mastery probability accordingly. Additionally, Duolingo's success with **A/B testing algorithmic changes** (e.g. switching to HLR improved user retention and presumably learning ¹⁸) highlights the importance of empirically validating any proficiency estimation method.
- **Coursera and edX (MOOCs):** Large MOOC platforms primarily measure achievement via **assignment grades and course completion**. They typically present *progress as a percentage of course completed* (which is mostly content consumption plus passing quizzes). However, in the background, MOOC research teams have extensively studied how **engagement metrics predict learning outcomes**. Examples:

- Students who watch most lecture videos and revisit difficult segments tend to perform better on assessments than those who skim or skip—video analytics research confirms **rewatching and slower, thorough viewing is associated with higher test scores**².
- Forum participation and reading supplemental materials are positive indicators of persistence and success. Some platforms generate an “engagement score” combining these activities to flag at-risk learners.
- No known MOOC gives a learner a *numeric proficiency score* purely from content views; instead, these metrics are used by instructors or platform algorithms to intervene (e.g. recommend content or send a warning if disengaged). This suggests our framework might use content-interaction data more for *formative assessment* and guidance rather than high-stakes certification. Nonetheless, the MOOC experience underlines that **early content engagement can predict final performance**⁴, meaning initial content interactions in a course could update a student’s estimated proficiency and allow tailoring of subsequent activities.
- **Brilliant.org and Codecademy:** These newer platforms emphasize **interactive learning by doing**, often embedding practice into content. Brilliant offers rich problem explanations and interactive visualizations. It doesn’t explicitly score “proficiency” per concept, but it does track which lessons a learner has completed and whether they solved the challenge. Codecademy similarly tracks progress through exercises (which are essentially coding tasks, i.e. applied questions). Both illustrate the trend of **blending content with immediate practice** – which is one way to get proficiency evidence. When pure content is consumed (e.g. reading a concept article on Brilliant), the platforms usually follow up with a quiz or have the learner apply it in the next step, thus indirectly measuring gain. For Nexus, if possible, inserting low-stakes questions or reflective prompts during content could greatly enhance measurement accuracy (though strictly speaking, that goes beyond “non-question” interactions).
- **Pearson’s Systems (e.g. MyLab, Revel):** Pearson’s own digital products often integrate readings with assessments. Revel, for instance, is an e-textbook platform that *intersperses quizzes in the reading*. The proficiency framework we design could draw on this approach by treating reading and interactive checkpoints together. Moreover, Pearson’s research (like the game log analysis with SPRING) shows they value **“invisible assessment”** – collecting evidence without interrupting learning⁹. Our framework aligns with that philosophy: use every interaction as potential evidence, reducing the need for separate tests. Pearson also has standardized proficiency scales (e.g. Global English Proficiency levels); while those are exam-based¹⁹, similar logic (a probabilistic interpretation of proficiency) can be applied to incremental content-based evidence.

Recommendations: (Detailed implementation guidance is provided in the main report, but key recommendations are highlighted below.)

- **Adopt a Hybrid Modeling Approach:** We recommend a **two-tier framework** – a base rule-driven model supplemented by data-driven refinement:
- **Rule-Driven Layer:** Immediately implement a **credit system for content interactions**. Assign provisional proficiency increments when students engage with content meaningfully. For example, finishing a tutorial video could award a small mastery boost to the related skills. Use domain expert input to set initial weights (e.g., watching a video might count as achieving 0.3 of a corresponding quiz question’s credit, reading a section = 0.2, completing an interactive demo = 0.5, etc., reflecting relative learning efficacy). Incorporate simple conditions to ensure quality: e.g., *only award points if*

video watched ≥90% or if scrolling behavior indicates entire article read. This provides an immediate, transparent way to value content engagement and can be communicated to learners (to incentivize using learning resources).

- **Data-Driven Layer:** In parallel, collect data to **train predictive models** and tune the rule layer. Use A/B tests and historical analysis to answer questions like: *Do students who receive content credit actually perform better later?; Which behaviors best predict subsequent question success?* For instance, apply regression analysis to see how time spent vs. quiz score improvements correlate. Over time, refine the weights or formulas in the rule-driven layer to improve accuracy. Eventually, a machine-learned model (e.g. a logistic regression or a simple neural network) could take multiple inputs (time, pauses, repeats, etc.) and output an estimated proficiency gain more optimally than static weights.
- **Bayesian Knowledge Tracing Integration:** Leverage existing question-based knowledge models in Nexus AI by feeding content interactions as “soft evidence.” In practice, if Nexus AI already uses a mastery probability per objective (as many intelligent tutoring systems do), introduce an update when content is consumed: increase the mastery probability by a small amount as per the earlier Bayesian formula. Maintain **separate confidence levels** for content-derived mastery vs. assessment-derived mastery – because content learning is unconfirmed until tested. One implementation is to give content interactions an effect on proficiency that decays over time *unless confirmed by a later correct answer*. This prevents the system from becoming over-confident from content alone. For example, reading about algebra might temporarily boost the model’s belief in the student’s algebra skill, but if the student never demonstrates it in a problem, the model could slowly relax that boost (modeling forgetting or uncertainty).
- **Continuous Calibration and Validation:** Establish a protocol to regularly validate the proficiency estimates against real performance. This includes comparing predicted vs. actual quiz results and running controlled studies when possible. For instance, if the model says a student mastered topic X from reading, give them a practice question on X – if many such students miss the question, the model needs adjusting (content was over-valued). Use these findings to calibrate the gain parameters ($\$L_c\$$ values in the Bayesian update or weights in ML model). As Khan Academy’s within-student study demonstrated, **increasing usage caused measurable test score gains** ¹, so we can back-infer reasonable credit: e.g., if doubling practice time yields +0.1 SD on an external test, the framework could distribute that expected gain across the content pieces completed.
- **Generalize Across Subjects and Ages:** The framework should remain **content-agnostic**, treating interactions with a science simulation or a history reading through the same lens of *time, engagement, and inferred understanding*. However, allow for **domain-specific tuning**:
 - In math/quantitative subjects, practice problems have especially high value (testing effect), so content reading might get relatively lower weight. In contrast, for humanities or languages, extensive reading or listening might itself build proficiency substantially (e.g. vocabulary growth from reading). Incorporate such differences by analyzing learning gain studies in each subject (for example, language learning benefits from immersive listening/reading, so those interactions should be weighted more heavily for language proficiency).
 - For younger learners (K-6), attention spans are shorter and reading skills developing, so interactive and visual content might be more effective. The framework can emphasize interactive content engagement (like virtual experiments, educational games) as evidence of learning for kids, whereas for adults, reading academic text might be effective. Ensure the model considers age-appropriate engagement patterns (younger students might need shorter content and more frequent assessment).

- Regardless of age or subject, focus on **foundational principles of learning**: active engagement, repetition, and feedback. Our framework gives credit for engagement, encourages repetition (if a student re-reads or replays content, that could indicate persistence and yield additional credit albeit with diminishing returns), and ultimately relies on occasional feedback (questions or external assessments) to stay on track.
- **User Feedback and Transparency:** Provide learners and educators with **clear insights** from the proficiency model. For example, show a dashboard that indicates: "You have spent 3 hours on Chapter 1 materials – estimated mastery 70% (confidence: medium). To verify and solidify this, you should now attempt the practice quiz." This way, the content-driven proficiency estimate is not a black box. It's accompanied by guidance ("we think you learned something, now confirm it by doing something active"). This philosophy follows the idea of using content interactions for *formative assessment* – it guides next steps rather than purely labeling the student. Teachers could also receive analytics highlighting students who engage a lot with content but still perform poorly on questions (flagging possible ineffective study habits), or vice versa those who don't engage and struggle (flagging low effort).

Conclusion of Executive Summary: Non-question interactions are valuable sources of evidence about learning, but they must be interpreted with care. By combining **educational theory** (how people learn from content) with **data analytics** (patterns that indicate real engagement and understanding), Pearson's Nexus AI can create a robust proficiency evaluation framework. The recommended hybrid approach will yield a system that is *broadly applicable across ages and subjects*, leveraging best practices from industry leaders like Khan Academy and Duolingo, while advancing beyond them by directly quantifying learning from rich content engagement. This will position Pearson's platform to provide personalized, accurate feedback to learners even in the absence of traditional assessments.

Introduction

Measuring student proficiency gains is traditionally associated with quizzes, tests, and assignments – explicit questions that assess knowledge. However, modern e-learning environments offer a wealth of **non-question learning activities**: students read interactive textbooks, watch lecture videos, explore simulations, write code, engage in discussions, and more. Pearson's Nexus AI, with its extensive content library, aims to harness these interactions to evaluate learning progress *in real time*. The goal is a **flexible proficiency evaluation framework** that remains valid without relying solely on test questions, thus enabling "stealth" assessment woven into learning activities.

Why focus on non-question interactions? There are several compelling reasons: - **Reducing Assessment Load:** Excessive testing can detract from learning time and increase anxiety ²⁰ ₈. By extracting assessment signals from learning activities themselves, we **"blur the line" between learning and evaluation** ⁹, aligning with educational best practices that emphasize continuous feedback over high-stakes exams. - **Continuous Personalization:** If we can estimate proficiency from content engagement, the system can adapt instruction on-the-fly. For example, if a student appears to grasp a concept just by reading an article (demonstrated by their engagement behaviors), we might fast-track them to more advanced material, whereas if another student speeds through a video with minimal interaction, the system might suggest additional practice on that topic. - **Inclusivity of Learning Styles:** Not all learning is demonstrated by test-taking. Some students learn deeply through reading or observing. A comprehensive

framework acknowledges **multiple forms of evidence** of learning, giving credit to students who are actively learning even when not quizzed at every turn.

Scope: This report covers **broad educational levels** (K-12, higher education, and adult learning), since the framework should adapt to different learner demographics. It remains **subject-neutral** – focusing on general principles and methods that apply across STEM, humanities, and language learning – while noting subject-specific nuances where relevant. We draw on case studies from major learning platforms (Khan Academy, Coursera/edX, Duolingo, etc.) to illustrate current practices and successful strategies. The report is structured as follows: - Background on prior research into learning from engagement and the theoretical basis for inferring knowledge from behavior. - Detailed exploration of methodologies to estimate proficiency gains from content interactions, with sub-sections on heuristic methods, analytical models, and machine learning approaches. - Platform case studies showing how industry leaders implement (or limit) these ideas in practice. - A proposed integrated framework for Pearson Nexus AI, including practical implementation steps, formulae for calculation, and system architecture considerations. - Validation strategies to ensure the framework's estimates are accurate and fair. - Supporting materials, including an annotated bibliography of sources used.

Throughout the report, **citations** to research literature and industry data are provided to justify and elaborate each point. These sources range from peer-reviewed studies in educational data mining to whitepapers and blog posts by ed-tech companies sharing insights from their platforms. This mix provides a balanced perspective combining rigorous academic findings with real-world application at scale.

Background and Theoretical Foundation

Inferring learning from behavior is a concept rooted in educational psychology and assessment theory. Two key paradigms frame our approach:

- **Learning Analytics & Engagement Theory:** In the last decade, learning analytics has emerged to make sense of big data from educational platforms. At its core is the idea that *every interaction a student has with content potentially signals something about their learning*. Researchers have identified various **engagement indicators** that correlate with learning gains or course success. For example, in online courses, metrics like **video interaction count, forum posts, and assignment views** have been used to predict final exam performance or course completion ^{5 2}. The underlying theory is that students who actively engage (spending time, revisiting material, etc.) are processing the content more deeply, which in turn leads to better understanding (consistent with *constructivist learning theory* – knowledge is built through active engagement). However, engagement metrics often show *weak to moderate correlations* individually ¹², meaning we should use them in combination and understand their limitations. Importantly, engagement can be superficial – hence the need to focus on **quality of engagement** (e.g. *self-explanation, reflection, deliberate practice* if such data can be captured).
- **Evidence-Centered Design (ECD):** ECD is an assessment design framework that guides how to use observable behaviors as evidence for unobservable competencies ⁸. In ECD, we define:
 - A **Competency Model:** the set of skills, knowledge, or abilities we care about (e.g. mastery of algebraic equations, understanding of WWII history).

- An **Evidence Model**: the behaviors or performances that indicate levels of those competencies (e.g. successfully solving an equation indicates algebra mastery; but what indicates it in non-question activities? Possibly spending time reading examples, or using the correct formulas in a simulation).
- A **Task Model**: the activities we present to learners to elicit that evidence (in our context, tasks are not just test questions, but also reading tasks, videos, projects, etc.).

Traditional tests have a straightforward evidence model: a correct answer is evidence of knowledge. For content interactions, we must craft an evidence model that says, for instance, “*viewing a video (task) and exhibiting certain behaviors (evidence) supports the claim that the student has learned X (competency)*”. This is challenging, but research like Shute’s **stealth assessment** in games provides a template: they built Bayesian networks with probabilities linking in-game behaviors to competency levels ¹⁰. That required expert input, but as mentioned, newer approaches aim to learn these links from data.

Cognitive Considerations: It’s useful to briefly review what cognitive science says about learning from non-question activities: - **The Testing Effect:** One reason inferring proficiency from passive study is tricky is the well-documented finding that *retrieval practice (testing) produces better retention than passive review*. Students often feel like they learn from reading or watching (familiarity increases), but without attempting retrieval or application, that learning can be shallow. Therefore, our framework should treat evidence from passive interactions as **lower-confidence evidence**. It’s not that no learning occurs – on the contrary, foundational knowledge often comes from initial exposure (you can’t solve a problem you’ve never seen content for). But we must be cautious not to equate “completed reading” with “mastered content” in a one-to-one way. - **Active Learning Strategies:** If we can detect *active learning strategies* during content consumption, those should boost the evidence value. For instance, research on reading shows that strategies like self-explanation, note-taking, or re-reading difficult passages improve understanding. In a digital system, proxies for this might be: highlighting text, writing a comment or note, clicking on glossary links, or replaying a segment of a video. These behaviors indicate the learner is mentally engaged and not just passively scrolling. Our evidence model should weigh such actions more heavily than simple completion. The study of handwriting logs in math noted that **writing speed and sustained focus time had a (weak but significant) correlation with proficiency** ¹² – suggesting that even how a student writes out notes or problems can reflect their cognitive processing. While Pearson’s platform may not capture handwriting (unless on a tablet), the principle stands: richer interaction data (beyond clicks) can give clues to cognitive engagement.

- **Cognitive Load:** The difficulty and novelty of content affect how much a student learns from it. If content is too easy (low cognitive load), a student might breeze through without gaining new skills; if it’s too hard (excessive cognitive load), the student might disengage or fail to comprehend. Recent research proposes **estimating cognitive load from user behavior and even physiology** (e.g. gaze patterns, heart rate) to better interpret learning events ²¹. For example, a moderate amount of struggle (indicated by re-reading a paragraph or replaying a video segment) could signal productive effort, whereas extremely long pause times or erratic navigation might signal confusion or overload. In our framework, if we have access to such data (perhaps in future via sensors or at least via patterns like repeated clicks on “rewind”), we could adjust the proficiency credit: *optimal engagement yields max learning gain, whereas signs of either disengagement or extreme struggle yield less gain*. This concept is advanced, but worth noting for future iterations of the model. A 2025 study integrating knowledge tracing with cognitive load found that balancing challenge with cognitive capacity led to better learning paths ¹⁴, underscoring that *how* the student engages matters, not just how long.

Summary: The theoretical foundation calls for a nuanced approach where: - We gather as many relevant behavioral signals as possible during content learning. - We assign evidence value to these signals based on a combination of prior research, expert judgment, and empirical calibration. - We update proficiency estimates in a probabilistic way, reflecting the uncertainty inherent in non-assessive evidence. - We continuously validate these estimates against real performance data (closing the loop so the model “learns” how to interpret behaviors more accurately over time).

In the next section, we dive into specific methodologies, building on these principles to outline concrete models and formulas for proficiency gain estimation.

Methodologies for Estimating Proficiency Gains from Content Interactions

In this section, we present and evaluate several methodologies, ranging from simple heuristics to sophisticated models. We also rank these approaches in terms of their **suitability** for Pearson’s needs, considering factors like accuracy, data requirements, implementation complexity, and explainability. Each subsection describes the approach, provides examples (with references to research or platform implementations), and discusses pros/cons.

1. Heuristic Point Allocation Systems

Description: A heuristic point system is the most straightforward approach: assign **points or scores for completion of learning activities**. For example, a student might earn 10 points for reading a section, 5 points for watching a video, etc. These points accumulate to indicate overall progress or proficiency. Often, points can be weighted by estimated difficulty or importance of the content.

Example Implementations: Many platforms use point systems for motivation. **Duolingo’s XP** is a prime example – users get XP for completing lessons (which involve both content and questions). While XP is primarily a gamification mechanic, it loosely correlates with effort spent. Another example is **Khan Academy’s energy points** (which historically were awarded for watching videos and practicing problems), though Khan has moved more toward mastery indicators now. **Codecademy** gives points for completing exercises, and **Udemy** shows a percentage completion for video lectures.

Formula: A simple scheme might be: - Each content item i has a base point value P_i (e.g. 100 points for a chapter, 50 for a short video). - If a student completes the item, they gain P_i . Partial completion could grant partial points (e.g. watched 50% of video = 50% of points). - The student’s proficiency in a topic could be the sum of points of items completed in that topic, possibly normalized to a percentage.

For instance, if Topic A has 5 content pieces (total 500 points) and the student earned 400 by thoroughly engaging with 4 of them, one might say they are 80% proficient in Topic A (pending any assessment confirmation).

Pros: - **Easy to implement and understand:** Both developers and learners grasp the idea of earning points. It provides instant feedback (“I finished this video, I got 10 points!”). - **Encourages engagement:** The system inherently rewards doing more, which can drive students to explore content. This is especially

effective for younger learners (points, levels, badges). - **No data dependency to start:** We can initialize this without historical data – expert judgment can set point values initially.

Cons: - **Low accuracy in measuring actual learning:** Earning points doesn't necessarily mean competence. A student could rush through content, collect points, but not internalize much. Without checks, you risk cases where a high “score” doesn't translate to high proficiency. - **One-size-fits-all weighting:** It might not distinguish between superficial and deep engagement. For example, simply marking a video as watched gives full points even if the student was zoned out. - **Inflation and scaling issues:** How do points translate to proficiency meaningfully? 1000 points means nothing by itself unless calibrated. Also, if content differs in difficulty, equal points might misrepresent their learning value.

Mitigations: To improve a point system, we can: - Tie points to **behavioral thresholds** (only award if the user met certain criteria, e.g., spent at least X minutes, or answered an embedded question correctly if available). - Use points as one component of a composite proficiency metric, rather than the sole indicator. For example, points indicate exposure while quizzes indicate demonstrated mastery – combine both for a final skill estimate. - Periodically **reset or rescale** points when course structure changes or when moving to advanced levels, so the meaning of “points” remains interpretable (some systems present them just as progress bar percentages per unit).

Assessment: As a standalone proficiency estimator, heuristic points are *ranked lower in reliability*. We recommend using this approach primarily as an **engagement tracker** and motivational tool, not as a high-stakes measure of skill. However, it forms a useful **base layer** – easy to implement and provides initial data to refine weights.

(Ranking: Good for quick deployment and motivation; poor for precision. Should be combined with other methods for actual proficiency evaluation.)

2. Time-on-Task and Engagement Quality Metrics

Description: Time-on-task is a classic metric: the duration a student spends on a learning activity. Numerous studies have attempted to correlate time spent with learning gains, with the intuitive notion that more time = more opportunity to learn. In practice, time alone is noisy, so modern approaches incorporate **engagement quality metrics** alongside raw time. These include: - **Completion percentage:** Did the student finish the video or read all pages of the text? - **Active time ratio:** How much of the time was the student actively interacting (scrolling, clicking, pausing) versus idle? - **Revisits:** Did the student come back to this content multiple times? Re-reading or re-watching may indicate reviewing for clarity (positive) or possibly confusion (could be positive up to a point, then negative). - **Pace indicators:** For videos, playback speed changes; for text, scrolling speed. Extremely fast completion might indicate skimming without understanding.

Example Implementations: - **MOOC Early Prediction Models:** A study by Yadav et al. (2020) (hypothetical reference) defined *attendance rate* and *utilization rate* for videos – essentially, how much of the available video content a student watched – and found they could predict course success with these metrics even after week 1 ⁴. High attendance/utilization early on signaled commitment and likely learning. - **Learning Management Systems (LMS):** Many LMS dashboards for instructors show which resources each student opened and for how long. While not directly turned into a “proficiency score,” instructors often infer understanding issues from this (e.g., a student spent unusually little time on a difficult reading might be

skipping it, or conversely, spent too much time might be struggling). - **Adaptive textbooks (e.g., Pearson's Revel):** They track reading time and give students feedback like "You spent 30 min on Chapter 3." Some even highlight if that's below class average, nudging more time. The assumption is minimum time is needed for adequate comprehension.

Formula: Building on the formula introduced in the Executive Summary for proficiency increment from content:

$$\Delta\text{Proficiency} = w_c \times f(\text{engagement metrics})$$

Where $f(\text{engagement})$ might be a composite score from 0 to 1 representing how fully the student engaged. For example: - $f(\text{engagement}) = \text{completionRate} \times \min(1, \frac{\text{timeSpent}}{\text{expectedTime}})$. - Here, completionRate is 1 if content finished, else a fraction. expectedTime could be a predetermined typical time to read/watch. The min ensures we cap credit at 100% even if they linger longer (to avoid gaming by just staying on a page). - We could multiply more factors: if r is number of replays (for a video) or page turns back (for text), perhaps include a factor like $(1 + \alpha r)$ up to some limit, under the reasoning that some review is beneficial (factor >1), but too many might indicate confusion (could invert effect if exceeding certain threshold).

Evidence from Research: - A comprehensive study on video clickstream data found that **total number of interactions (clicks)** had positive correlation with test performance ²¹. This suggests that simply *doing stuff* (clicking pause, rewind, etc.) often reflects engagement. Notably, the *number of pauses and slow rewinds* were among the strongest predictors of higher test scores ²², whereas frequent skipping ahead (fast-forward) was associated with lower performance. This implies that *thorough processing of content yields measurable learning benefit*, which our model should reward. - Conversely, the study noted these correlations, while significant, were not extremely high (correlation coefficients in the 0.3-0.5 range typically). That's expected, as many other factors influence learning. It reinforces that time/engagement metrics should be used in aggregate and with calibration.

Pros: - **Data is continuously available:** Time and basic interaction logs are collected by all digital learning systems. No special setup or instrumentation beyond what already exists in Nexus AI. - **Granular adjustments:** We can give partial credit. For example, if a student only watched 50% of a video, we can assign half the potential proficiency gain. This reflects the intuition that incomplete use of content yields partial learning. - **Adaptability:** We can adjust expectations for different content. Short tasks vs. long tasks, we normalize time by expectedTime as above. We can personalize expectedTime too – if a student reads slower (based on past behavior), spending less time might still mean they read thoroughly. A sophisticated model might learn each student's typical reading speed to better interpret time (this borders on ML approach). - **Qualitative interpretation:** Teachers and learners can understand a message like "We credit you more because you spent adequate time and didn't skip around". It aligns with common-sense study advice (take your time, don't skip material).

Cons: - **Potential for gaming or misinterpretation:** A student might realize time spent is tracked and simply leave a video playing to get credit, or scroll slowly through text without reading just to appear engaged. Without secondary checks, the model could be fooled by such behavior. We would need to implement some detection (e.g., if video is in a background tab or if there's no mouse/scroll activity for long periods, discount the time). - **Variability in needed time:** Some students learn quickly and don't need as much time; others may spend a lot of time but still be confused. So equal time doesn't equate to equal

learning for everyone. That's why combining multiple metrics (time + completion + replays etc.) is necessary to refine the signal. - **No direct content understanding measure:** Time on page doesn't tell us if the student grasped the concept. Two students could spend 10 minutes, one deeply processing, another daydreaming. We mitigate this by looking for signs of active engagement (clicks, highlights, etc.), but it's still indirect.

Use in Framework: Time-on-task with engagement metrics is a **foundational component** of our recommended approach. It should be used to modulate the **initial proficiency boost given by content interactions**: - Define thresholds (e.g., must spend at least X minutes or Y% of video to count as "completed"). - Reward additional engagement like second views (perhaps with diminishing returns: the first view gives 100% of content's weight, a second full view might give an extra 50%, etc., acknowledging review helps but not as much as initial learning). - Integrate with the knowledge model: the L_c "learning probability" for content could be made a function of engagement quality. For example, if normally reading a section has $L_c = 0.3$ chance to learn the skill, a highly engaged read might effectively raise that to 0.4, whereas a perfunctory skim might be 0.1.

(Ranking: High importance in framework – provides continuous signal. Needs careful design to prevent superficial engagement from inflating proficiency. Overall, an essential middle-ground approach between naive points and complex models.)

3. Embedded Knowledge Tracing and Probabilistic Models

Description: This approach integrates content interactions into a formal **knowledge tracing model** or uses Bayesian networks / IRT (Item Response Theory) to update skill proficiencies. Essentially, we extend methods traditionally used for question-answer data to also handle "learning event" data.

Bayesian Knowledge Tracing (BKT) Extension: In standard BKT (used in systems like Cognitive Tutors), each skill has a binary state (learned or not learned) with an associated probability. Each practice question is an observation that updates this probability via Bayes' rule, considering slip/guess probabilities. BKT also includes a learning parameter – the chance the skill is learned *during* a practice opportunity even if the answer was wrong (e.g., by receiving feedback or by the attempt itself).

For content, we can analogously say: each content interaction is a learning opportunity with a certain probability of causing the skill to transition from unlearned to learned. If content is consumed, we update the skill's probability: - If the skill was not mastered before, now there's a chance it is mastered (increase the probability by some factor). - If it was mastered, it stays mastered (or possibly even reinforces it, though BKT usually assumes once learned, always learned – forgetting can be modeled separately).

This is effectively the formula mentioned earlier:

$$P_{new}(\text{mastery}) = P_{old}(\text{mastery}) + (1 - P_{old}) \times L_c,$$

where L_c is the learning probability attached to content piece c .

Item Response Theory (IRT) Analogy: In IRT, a student's ability and an item's difficulty determine probability of a correct response. While IRT is for questions, one might imagine treating a content piece as an "item" that the student can "master" or not. Since there's no wrong/right, we can't directly fit IRT, but we

could invert the logic: if a student later answers questions on the material correctly, that "validates" that the content was effectively learned from. If not, perhaps it wasn't. This is more of a calibration mechanism than a direct model, so BKT is more straightforward in implementation.

Bayesian Networks (multi-skill): Some content covers multiple skills. Bayesian network models (or Dynamic Bayesian Networks for sequences) can represent multiple knowledge components and evidence relationships. For example, reading a complex article might increment two separate skill probabilities. This requires content-skill tagging and possibly different weights per skill (the article might primarily teach skill A with $L_c = 0.5$ but also tangentially benefit skill B with $L_c = 0.2$). Setting those values would initially rely on expert guess or prior studies (or how integral that content is to each skill, maybe based on curriculum mapping).

Pros: - **Principled, interpretable model:** Unlike a raw ML approach, a Bayesian knowledge model is easier to interpret and explain to stakeholders. Each skill's probability has meaning, and updates are mathematically principled. - **Combines with question data naturally:** This approach doesn't live in a silo – it can seamlessly merge with question-answer data. For instance, after a series of content interactions, a quiz response can be used to update the same probability, thus validating or adjusting the earlier inferences. - **Allows uncertainty modeling:** If a student only read content and never answered a question, the skill probability might be, say, 0.4 – indicating we're not sure of mastery yet. The model naturally expresses that uncertainty. It can also accumulate evidence: reading two different sources on the same skill might push probability higher than just one, approaching a mastery threshold.

Cons: - **Requires parameter tuning:** What should L_c be for a given content piece? If set too high, the model will overestimate learning (false mastery flags); if too low, it undervalues content. These might differ by content difficulty and by student prior knowledge. Tuning these parameters likely requires data (comparing model predictions to actual outcomes) or at least iterative refinement by experts. - **Independence assumptions:** BKT assumes each learning opportunity is independent given the current state, and it doesn't directly use the rich sequence info beyond the probability state. If a student flounders through a video (which we could detect via behaviors), standard BKT wouldn't know that – it would just see "content completed" and apply the same L_c . We'd need to extend BKT or use a more flexible model to incorporate engagement level (one could have multiple versions of each content event – e.g. "content thoroughly engaged" vs "content skimmed" – with different L_c values). - **Forgetting and revision:** Proficiency isn't static; a student might "learn" a concept from content but forget it if not practiced. BKT can incorporate a forgetting factor or we re-assess via periodic questions. It's wise to include a *decay over time* for probabilities in the absence of reinforcement, which complicates the model but aligns with reality (and is something Duolingo's model handles via half-life decay).

Use in Framework: We consider this a **top-ranked approach for accuracy**, especially when combined with periodic assessments. The recommended framework will incorporate a Bayesian updating mechanism. Concretely: - Pearson Nexus could maintain a vector of skill mastery probabilities for each student. Initially, based on maybe a placement test or assumed baseline. - Each content interaction triggers an update: for each skill tagged in that content, adjust its probability upward by a factor relative to engagement quality (as discussed, if engagement was poor, perhaps we treat it as a very low L_c or even skip updating). - The system can define mastery thresholds (e.g. if $P > 0.8$, consider skill mastered) but also require some confirmed evidence (like at least one quiz correct) to mark it officially mastered – a policy decision to ensure reliability. - An **example:** Student's probability for "Photosynthesis concept" is 0.2. They watch an interactive animation about photosynthesis with full engagement. That content has $L_c = 0.5$ based on prior

calibration (i.e., we believe a thorough viewing gives them a 50% chance to learn it). The new probability becomes $0.2 + (1 - 0.2) * 0.5 = 0.6$ (a sizable jump, reflecting the significant instructional content). We're still uncertain if they truly got it (60%), so perhaps we suggest a practice question. If they then answer correctly, the model would spike them higher (maybe to ~0.85, considering the possibility of guessing). If they answer incorrectly, the model would drop the probability (maybe back down to 0.3, indicating the content alone didn't suffice).

(Ranking: High. This method is core to a precise, individualized tracking of proficiency. It leverages well-established techniques in intelligent tutoring systems. The only reason not to use it exclusively is the effort to set it up and tune it; however, given Pearson's resources and the importance of accuracy, this is a recommended backbone. We will calibrate it with data and possibly simplify for initial deployment by grouping some skills or using IRT-like simplifications if needed.)

4. Machine Learning and Predictive Analytics Models

Description: This category encompasses using data-driven models to learn the mapping from content interaction patterns to proficiency outcomes, rather than relying on human-designed formulas or simple probabilistic updates. Techniques include: - **Regression models:** e.g. linear or logistic regression taking features like total time, number of content pieces completed, etc., to predict a continuous proficiency score or probability of passing a test. - **Classification models:** predicting a categorical outcome (pass/fail, high mastery vs low mastery) based on engagement metrics. - **Sequence models:** e.g. RNNs or Transformers that take the entire sequence of a student's activities as input and predict their knowledge state or next performance. (DKT mentioned earlier is one such model, though it typically uses question attempts, it could be extended to include other event types in the sequence). - **Recommender-system style matrix factorization:** treat proficiency as a latent trait and try to factorize an interaction matrix (this is more abstract, but one could imagine content pieces as "items" and whether the student engaged or not as implicit feedback, trying to infer an ability parameter).

Example Implementations & Research: - The **SPRING model by Pearson** (Gonzalez-Brenes et al. 2016) essentially automated evidence model creation from game logs ²³ ²⁴. They used a combination of clustering and classification to predict post-test performance from in-game behavior, achieving a decent accuracy (Correlation 0.55 with test scores) ⁶. This is a proof of concept that ML can extract signal from complex clickstream data beyond what a human might easily identify. - **MOOC dropout and grade prediction:** Many studies train classifiers on week-by-week behavior to predict who will earn a certificate. They input counts of videos watched, % of assignments done, forum words written, etc. Results often show 80-90% accuracy in identifying top performers vs dropouts by mid-course. One study cited earlier found using the first week's video stats predicted passing with high effectiveness ⁴. - **Duolingo's half-life regression (HLR):** While HLR is a custom model blending theory and ML, it essentially performs a regression on features of practice history to predict memory strength ¹⁶ ²⁵. They fit this model using millions of data points, demonstrating that large-scale data can inform the model parameters better than guesswork. HLR improved prediction of recall (word strength) by nearly 50% over older methods ²⁶. This indicates that ML-based parameter tuning, even within a theoretical model, can significantly boost accuracy.

Inputs (Features): The richness of ML is that we can incorporate **many features**: - Simple counts: # of content items viewed, # of repeats, total time, etc. - Sequence features: e.g. did the student use an optimal study pattern (spaced reviews) or cram? - Content difficulty: we can include metadata, like reading level of text or length of video, to let the model learn how those interact with mastery. - Learner behaviors: e.g.

whether a student tends to take notes (if such a feature exists in the platform), or their average pace of learning historically (some students consistently under-engage with content). - Demographics or prior knowledge as features: possibly age, grade level, or entry test scores, so the model can compare engagement relative to expected.

Output: Could be a predicted test score, or probability of mastery for each skill. We might train separate models per skill (if enough data per skill) or one overall model that outputs something like “probability student will answer a random question from this topic correctly”.

Pros: - **Potentially high accuracy:** ML can capture complex nonlinear relationships and interactions between variables. For example, maybe time only helps if above a threshold, or maybe videos help more for visual learners (if we have any proxy for that). A human-crafted model might miss these nuances that an algorithm could pick up. - **Automatic discovery of patterns:** ML might find unexpected indicators. For instance, a model could learn that *frequent pausing in videos is a strong positive sign* (consistent with what Yürüm et al. found ²) whereas a certain pattern like “skips all videos but spends long on reading” might correlate with certain outcomes. These insights can refine our understanding. - **Adaptive to data from Pearson's platform:** Every platform has unique usage patterns. A model trained on Nexus AI data will inherently tailor to Pearson's content and users, potentially outperforming generic models. It can also be updated periodically as more data flows in (e.g., retrain each semester). - **Incorporating multivariate outcomes:** If Pearson has multiple ways to gauge proficiency (like periodic diagnostics or external exam scores for some students), ML can use those as ground truth to optimize against, potentially combining signals in an optimal way to predict those outcomes.

Cons: - **Data hungry:** Training effective models requires a lot of historical data where we know both the input (engagement logs) and output (true proficiency, typically measured by tests). For new courses or new content, we may lack enough data. Cold start is a problem – initially, we rely on theory until data grows. - **Opacity:** Many ML models, especially neural networks, are black boxes. It might be hard to explain *why* the model thinks a student has 80% proficiency. This can be a concern for teacher trust and for debugging. We might mitigate with simpler models (like decision trees or regressions) or use explainable AI techniques to interpret feature importance (e.g., SHAP values). - **Overfitting and bias:** If not carefully validated, models might latch onto spurious correlations. E.g., they might learn that “students who click help a lot have low scores” which could be true on average but if interpreted causally could discourage help-seeking which is not our intent. We must ensure the model is used for prediction, not to judge behaviors out of context. Also, ensure fairness: if some students (like non-native English speakers) naturally spend more time reading content, the model shouldn't misinterpret that as either always good or always bad without context. - **Maintenance:** As content changes (new courses, revised chapters), the model must be retrained. If features aren't carefully defined, changes in platform (like new interaction types) can break the model. This is a long-term maintenance cost and requires AI expertise on the team.

Use in Framework: We envisage ML models as a **validation and enhancement layer** rather than the primary real-time calculation (at least initially). The recommended path: - Start with the interpretable rule-based/Bayesian framework as a backbone (sections 1–3 above). - In the background, collect data on how well that framework's estimates align with actual performance data. Then, build ML models to analyze residuals and suggest improvements. - For example, a regression could be trained where the dependent variable is the student's eventual test score, and independent variables include the current framework's estimated proficiency plus various raw engagement metrics. If the regression finds that certain metrics have additional predictive power (meaning our framework under-utilized them), we can incorporate those

into the model explicitly. Or we might even replace some heuristic calculations with learned formulas from regression. - In time, once enough trust and understanding in the model is built, we could move to a scenario where a trained model directly outputs proficiency estimates in real-time. This might be feasible for well-trafficked courses where large data is available. For long-tail or new courses, we'd fall back to the robust Bayesian model with default parameters.

An example ML-driven enhancement: Perhaps the data shows that *watching at least 2 related videos* has a nonlinear effect – it's particularly good if the student does both, versus just one. Our initial model might not have a rule for that. A decision tree could discover "if *videos_watched >= 2* then *proficiency boost += X*". We can then encode that rule in the system.

(Ranking: Medium-High. In terms of ultimate effectiveness, this could be top-notch if executed well. But due to its prerequisites (data, expertise) and concerns of transparency, it should augment, not replace, the structured approaches. We rank it as a priority for research and iterative improvement of the system, but not the first thing to deploy on day one.)

5. Qualitative and Self-Reported Measures (Supporting Approach)

(We include this for completeness, though it's a somewhat orthogonal approach to the data-centric ones above.)

Description: This involves gathering either **self-assessments from students** or observational assessments from teachers about the student's proficiency gained from content study. While not a direct automated method, these can complement algorithmic approaches: - After a student engages with content, ask them a reflective question: "Do you feel you understood this material?" or "Rate your confidence in this topic now." Such metacognitive prompts can both reinforce their learning (by making them reflect) and provide a signal. Research shows students are reasonably good at judging if they completely didn't get something vs. if they feel comfortable, although self-assessments can be over- or under-confident. - Teachers (or the system) could prompt short summaries or concept maps from the student after reading, which can be evaluated (possibly by AI or rubric) to gauge understanding.

Use in platforms: Some e-learning platforms occasionally ask "Was this article helpful to you?" or "Quiz yourself: write the key takeaways." These are not widespread and often optional. However, Pearson might integrate something like a "Lesson checkpoint: write 1-2 sentences about what you learned" – this becomes qualitative data that either a teacher reviews or NLP techniques attempt to parse for correctness.

Pros: - **Direct insight into student thinking:** Especially open responses can reveal misunderstandings or correct but unassessed knowledge. - **Encourages active recall (which itself boosts learning):** Asking a student to summarize or answer a quick reflective question after content essentially turns it into a mini assessment (blurring our line, but still a non-traditional one). - **Calibration of student confidence:** Over time, comparing self-rated understanding with actual quiz performance can help calibrate the student's self-awareness and also the system's trust in self-reports.

Cons: - **Subjective and inconsistent:** Students might lie or misjudge. Some will always say "got it" even if they haven't (to move on or due to overconfidence); others might say "still unsure" even when they actually understood (perfectionists). - **Not automated at scale (unless using AI grading for summaries):** Getting value from open-ended responses is non-trivial. It could become a text analysis problem if we go that route (which loops back into ML/NLP). - **Intrusive to user experience:** If overused, asking students to constantly

self-evaluate or write summaries might annoy them or slow them down. It has to be used sparingly and meaningfully.

Recommendation: We do not rely on self-reports as a primary proficiency indicator, but we consider using them as a **supplementary signal** and for research purposes. For instance, we could log whether a student chose to re-watch a video or read an optional “related topic” – that itself is a form of self-regulation signal: the student felt they needed more, implying they hadn’t mastered it from the first content. So indirectly, student choices can be treated akin to self-assessment (if they rush ahead, they think they know it; if they seek more material, they likely felt unsure). Our ML models could incorporate such signals (did the student view additional resources or not).

(Ranking: Low for automated proficiency computation, but potentially useful when combined with other data. Could be part of future expansion focusing on metacognition.)

Having outlined these methodologies, **Table 1** (below) provides a side-by-side comparison, including criteria like implementation complexity, data requirement, immediacy of feedback, and proven effectiveness:

Approach	Accuracy	Ease of Implementation	Data Needed	Notes
Heuristic Points (XP)	Low (standalone)	Very Easy	None (theory-based)	Good for engagement, not reliable proficiency measure on its own.
Time & Engagement Metrics	Medium (with good design)	Easy-Moderate	Low (basic logs)	Needs quality filters; forms basis of many predictive patterns.
Bayesian Knowledge Tracing	Medium-High	Moderate (complex logic)	Moderate (needs tagging & tuning)	Provides individualized, interpretable skill estimates; tunable with data.
Machine Learning Model	High (if well-trained)	Hard (requires expertise)	High (large historical dataset)	Can discover complex patterns; risk of overfitting; use to refine other methods.
Self-Report/Qualitative	Variable	Easy (for self-report) / Hard (for analysis)	Low (for asking) / High (for analyzing)	Useful for reflection and supplementary data; not core automated metric.

Table 1: Comparison of proficiency estimation approaches from non-question data.

Our recommended framework will blend the first four approaches: using time and engagement as immediate signals, feeding into a Bayesian update model for proficiency, and continuously improving via machine learning on accumulated data.

Case Studies: How Learning Platforms Estimate Proficiency from Engagement

To ground the discussion, we examine how several well-known learning platforms incorporate non-assessment interactions into their learning evaluation (or why some deliberately don't).

Khan Academy

Model: Mastery-Based, but *assessment-centric*. Khan Academy's system is built around mastery of skills through practice exercises. A student's proficiency in a skill goes from "not started" to "familiar" to "proficient" to "mastered" based on answering questions correctly, especially in Mastery Challenges (cumulative quizzes). Non-question activities (videos, articles) do not directly raise the skill level ¹⁵. They are available resources, and teachers can assign them, but the student doesn't get mastery credit just for watching.

Rationale: Khan Academy's stance is that **watching a video doesn't prove mastery**, only solving problems does ¹⁵. This is pedagogically aligned with the testing effect. From a proficiency standpoint, Khan uses content engagement mainly to *help the student learn so they can do the exercises*. In the data they've published, students who spend more time using Khan (which usually means a mix of videos and practice) see higher learning gains externally ¹. But internally, their tracking of proficiency remains tied to problem performance.

Notable Features: - They do provide **analytics** to teachers that show time spent on videos, which can be an early indicator. A teacher might notice a student never watches the videos and struggles in practice, prompting an intervention ("maybe watch the video on this concept"). - Recently, Khan Academy has integrated an AI tutor (Khanmigo) which engages students in conversations about content. This introduces new interaction data. It's likely Khan will need to evaluate learning from dialogue (e.g., did the student ask good questions? Did they arrive at correct understanding in the conversation?). While no public info is available on how they measure this, it's plausible they'll use some heuristic (if the student's answers in the dialogue are correct after some hints, then mastery is improved).

Takeaway for Pearson: Khan's approach represents one end of the spectrum – *cautious and assessment-heavy*. It ensures high confidence in mastery but might under-utilize rich data (e.g., if a student diligently watches 10 videos before ever attempting a question, Khan's system knows nothing of their knowledge until the first question). For Nexus, we likely want a more continuous read on the student (to offer adaptation sooner), so incorporating content interactions is beneficial. Still, Khan's approach reminds us that any inferred proficiency from content should eventually be checked by application.

Coursera and edX (MOOCs)

Model: Progress and completion based, with **post-hoc analytics**. MOOCs typically give a final grade determined by assignments and exam performance (questions). Content viewing isn't part of the grade.

However, platforms track it and researchers analyze it. They often give instructors **dashboards** with engagement metrics: - e.g., **Video engagement**: what % of a video was watched by each student, how many attempted the embedded quizzes (if any), etc. - **Resource access**: which papers or pages were clicked, how long in forums.

During a course run, instructors might use this data to identify students who are less engaged (to send emails or adjust content). At scale, automated interventions have been trialed: like sending “nudges” to students who fall behind in watching videos (“Week 3 videos are now up, we noticed you haven’t finished Week 2’s videos yet...”).

Research and findings: - MOOC studies found strong **correlations between engagement and success**: students who complete most videos and assignments are exponentially more likely to pass (which is somewhat tautological since assignments determine passing). More interestingly, *early engagement is predictive*: someone active in the first two weeks is far likelier to stay and pass ⁴. - Fine-grained behaviors in videos have predictive value, as noted: a student who consistently pauses and rewinds might be more reflective and thorough, hence learning better ². MOOC platforms haven’t explicitly productized that insight (like there’s no “video engagement score” given to students), but it informs design (e.g., edX has experimented with interactive transcripts, allowing students to click back on certain terms – this could be seen as good engagement).

- **Open edX** (the open-source platform) has an “**engagement score**” plugin that was developed in research, which gives a single score factoring in videos watched, problems attempted, forum posts, etc. This was used to predict who might need help. It’s not widely known to end-users, more for research/instructors.

Takeaway: The MOOC approach is to treat content interaction as **ancillary data** to support student success, rather than a metric of proficiency. The end goal proficiency (getting a certificate or passing) is still judged by assignments and exams. For Nexus AI, if the use-case is to provide ongoing estimates of student proficiency (perhaps to guide them or inform instructors), relying solely on eventual tests defeats the purpose of early intervention. So, we lean more into using the engagement data directly. But the MOOC data suggests which signals are most useful to look at (video watch patterns, timely engagement with materials, etc.). Also, because MOOCs have very diverse adult learners, it shows that these patterns hold broadly, not just for K-12 – giving confidence our approach can be broad.

Duolingo

Model: Gamified practice with **adaptive review**. Duolingo’s primary measure of progress is XP (experience points) and **skill crowns** (levels achieved in each skill through practice exercises). All of those come from actually doing the interactive exercises (translating sentences, etc.). Pure content (like reading tips about grammar) is optional and not explicitly scored. However, Duolingo does maintain a hidden model of the user’s knowledge strength per word/skill – the Half-Life Regression (HLR) model we discussed. That model *only updates when you practice items* (get them right or wrong), since it’s fundamentally about memory decay after practice.

Interesting aspects: - Duolingo’s system will schedule content for you. If a student is just reading tips and not practicing, the app will eventually prompt practice because it doesn’t count reading as evidence. In a sense, Duolingo forces a question-based check for everything. - That said, Duolingo’s design intermixes

content and practice tightly (microlearning). A lesson might introduce a new word with a picture (content) and immediately ask you to use it. The immediate attempt is both learning and assessment. This design highlights that an optimal way to measure learning from content is to follow it quickly with an application opportunity, reinforcing the learning and giving the system data. For Nexus, where possible, embedding small questions or interactions in content (even as ungraded self-checks) could dramatically improve our proficiency estimates. If such embedding is not possible or desired (maybe to keep the experience pure reading flow), then we must rely on less direct signals, which as we know, are weaker.

- **Data analysis:** Duolingo has published that their HLR model improves prediction of which words users will forget ²⁶. Indirectly, that means they can estimate proficiency of a word at any time. But because it's reliant on practice history, it again underscores the point: they needed explicit practice data. Duolingo does not try to gauge if you learned a word just by hearing it once in a story without you then using it.

Takeaway: Duolingo's success is largely due to its *active learning loop*. For our context, if Nexus AI can encourage even minimal active responses during or after content, it will enhance the proficiency measurement greatly. But if we assume we truly have to evaluate proficiency from passive interactions alone, we should be more conservative (like Khan and Duolingo are) in our confidence. It might be that the best we do is say "the student likely improved their proficiency by reading this, but we need to verify with a question or task soon."

However, Duolingo's HLR gives us a template for dealing with knowledge decay and practice schedules. We might incorporate a similar decay factor: if a student's only evidence on a topic is that they read something 3 weeks ago, our framework should probably decrement their proficiency estimate slightly over time, until refreshed by either re-reading or a successful question.

Other Platforms (Brilliant, Codecademy, etc.)

- **Brilliant.org:** Focuses on STEM problem solving with interactive explorations. Brilliant's measure of progress is basically completing topics and solving challenge problems. If you just read their explanations and don't attempt problems, you won't "complete" the lesson. They often present a problem first, then explanation, then another problem. This suggests that they consider the explanatory content as assistive, not as evidence. They do track if you viewed the explanation (and you earn "points" on the site for participation), but skill proficiency is evidenced by solving the problems. In user profiles, Brilliant shows "topics completed" and "quizzes solved" more prominently than any content consumption.

Lesson: Structure content to include interactive checkpoints, or at least note that even highly regarded platforms lean on problems for assessment. But since Pearson's content library might include things like textbook chapters, we want to get as much info as we can from reading them – perhaps using analogies like "solved examples = evidence".

- **Codecademy:** Teaches programming by combining instructions with live coding exercises. Every exercise is essentially a task (write code to do X) which is automatically checked. So Codecademy can gauge proficiency by whether the student's code passes tests. Reading the narrative or hints is not tracked as learning progress except insofar as you eventually solve the tasks. This again is an active learning model.

- **SMART Sparrow / Adaptive Labs:** Some platforms that focus on simulations or virtual labs (e.g., in science education) use **embedded metrics**. For example, a virtual physics lab might track if the student discovered the right concept by how they manipulate variables. If they eventually set up the simulation to get the expected result, one infers they learned the principle. These are quite domain-specific assessments (almost like performance assessments). If Pearson has interactive simulations, we could set up criteria: e.g., “student successfully balanced the chemical equation in the simulation” as a binary success, thus proficiency credit. That crosses into actual problem-solving, though.

In summary, the trend is clear: platforms often *don't trust content consumption alone* to determine proficiency; they pair it with questions or tasks for confirmation. Those that do use content data (like research projects or low-stakes analytics) find it useful for prediction and early warning, but rarely as sole proof of mastery.

Positioning Pearson's Nexus: There is an opportunity for Pearson to innovate by more explicitly leveraging content engagement data in a systematic way. If we succeed, Nexus AI could provide a much richer picture of student learning in real time, especially in content-heavy courses (where traditionally you'd only know how the student is doing when they submit an assignment or take a test).

Our framework takes inspiration from these case studies by combining: - Khan's rigorous linking of mastery to evidence (we'll ultimately validate via questions). - MOOCs' use of engagement metrics to anticipate performance (we'll use similar metrics to adjust proficiency continuously). - Duolingo's idea of modeling memory and practice (we incorporate forgetting curves or practice schedules so content learning is time-sensitive). - Brilliant/Codecademy's approach of immediate practice (we recommend wherever possible to incorporate practice, even if minor, around content – but when not possible, we at least observe behaviors as pseudo evidence).

Next, we detail the proposed framework and how Pearson could implement and calibrate it, followed by how to validate its effectiveness.

Proposed Proficiency Evaluation Framework for Nexus AI

Based on the research and analysis above, we propose a comprehensive framework that integrates multiple approaches to estimate proficiency gains from non-question interactions. The framework can be conceptualized in the following components or stages:

1. **Content Tagging and Metadata:** - Ensure each content piece (video, article, simulation, etc.) is mapped to the relevant **learning objectives or skills**. This mapping could be one-to-one (one content corresponds to one skill) or one-to-many (a chapter covers several skills). Pearson's existing curricular metadata can be leveraged here. Granularity is key: the smaller the skill units, the easier to update specific proficiencies. For example, a 10-page chapter might map to 3 distinct skills – we'd track progress on each. - Assign initial **weight/importance** to each content piece regarding each skill. This could be a value of L_c (learning probability) as discussed, or simply a nominal “learning impact score.” For instance, a comprehensive tutorial video might have a higher impact score than a brief summary text. - Record expected engagement parameters: e.g., expected reading time, video length, difficulty level. These will be used to gauge whether a student's interaction was sufficient.

2. Student Engagement Monitoring: - As a student engages with content, the system logs detailed events: start time, end time, completion percentage, pauses, rewinds, scrolls, clicks on interactive elements, etc. Nexus should unify these logs across content types into a common event stream (possibly in Experience API (xAPI) format or similar, which is suited for capturing learning experiences in a standardized way). - Compute **engagement metrics in real-time**. For each content item, derive an engagement score (0 to 1) that encapsulates how thoroughly the student engaged. We defined examples earlier; to recap one possible formulation:

$$\text{engagementScore} = \text{completionFraction} \times \min\left(1, \frac{\text{timeSpent}}{\text{expectedTime}}\right) \times f(\text{interactionDepth})$$

where $f(\text{interactionDepth})$ might be a function that increases with meaningful interactions (pauses, notes) and decreases if the student skipped or rushed. - Determine **completion status**: Did the student complete the item meaningfully, partially, or not at all? This status could be used in gating (e.g., require completion before marking a unit done) and also in proficiency updates (only give full credit on substantial completion).

3. Instant Proficiency Update (Rule-Based): - When a student finishes engaging with a content item, the system immediately triggers a proficiency update for the associated skills. This uses a rule-based formula initially. For each skill s linked to content c :

$$P_{\text{new}}(s) = P_{\text{old}}(s) + (1 - P_{\text{old}}(s)) \times L_{c,s} \times \text{engagementScore}$$

Here $L_{c,s}$ is the content's learning probability for skill s (the maximum potential increase if the skill was unmastered and engagement was perfect). engagementScore is as computed. $P_{\text{old}}(s)$ is the prior mastery probability. - This formula ensures diminishing returns: if P_{old} is already high, the gain is smaller. It also ensures that no matter how high $L_{c,s}$ is, the skill probability never exceeds 1 or the update overshoots. - If multiple content pieces related to the same skill are studied, each provides an incremental update. (We might cap at some asymptote if necessary to prevent spurious inflation without assessment.) - If skills have prerequisites or are hierarchical, we might also propagate some effects – e.g., learning a foundational concept could partially raise advanced concept probabilities, or vice versa. This could be modeled in a Bayesian network if needed.

- **No immediate negative updates** are given, since consuming content can't directly reduce knowledge (except by forgetting over time, handled separately). However, if the system detects extremely poor engagement (e.g., student skipped through quickly), we might decide not to increase (or increase only a negligible amount, effectively treating it as no learning occurred).
- Log the update along with a rationale that can be shown to the student or teacher: "Read Section 3.2 – increased mastery of Skill X from 40% to 55% (partial understanding expected from reading)."

4. Confidence and Flags: - Alongside proficiency $P(s)$, maintain a **confidence level** or evidence count. For example, label whether this skill has been "confirmed" by a direct assessment. If a skill's high probability is based solely on content interactions, mark it as *unconfirmed*. This can be a simple flag or a confidence score that is lower in such cases. - Use these flags to prompt assessments: e.g., if a student has content-based

mastery $\geq 80\%$ on a skill but no practice question yet, the system can say “Time to verify your knowledge – try this problem.” This ensures that content learning is eventually vetted. - Alternatively, incorporate the lack of direct evidence in the probability itself: one could adjust $P(s)$ downwards until confirmed. But it might be cleaner to keep a separate measure: e.g., show mastery as “75% (unverified)”.

5. Periodic Assessment & Calibration: - When the student does answer questions or completes an assignment, the standard knowledge tracing update occurs. If they get it right, great – our content-inferred proficiency is validated (and likely gets a further bump if not at 100%). If they get it wrong, the framework must reconcile this: - Possibly lower $P(s)$ significantly, because this suggests the content learning didn’t solidify true mastery. - Mark that content-only evidence was not sufficient here, which could signal to adjust L_c values for that content in the future or prompt recommending additional content to the student. - Over many students, gather stats: e.g., among students who only learned via reading content and had $P(s) \sim 70\%$, what percentage answered a related question correctly on the first try? If it’s only 40%, then 70% was an overestimate – maybe we should lower L_c or adjust how we calculate that probability. This kind of **calibration analysis** will refine the model. It can be done offline by data scientists periodically and updates pushed to the rule base.

- Similarly, track if content consumption yields *no apparent gain*. If lots of students read a certain piece but their skill doesn’t improve (they still get questions wrong), that content might be low-quality or misaligned. This crosses into content efficacy analysis, an added benefit of having this data.

6. Machine Learning Enhancement (Long-term): - As data accumulates (e.g., a semester of students using the system), train ML models to predict assessment outcomes or final mastery from the full history of content interactions. Use this to adjust the weights: - Perhaps a model finds that “for skill Y, watching videos was twice as effective as reading text in terms of likelihood of getting quiz questions right.” Then we might set the L_c for videos on skill Y higher relative to text. - Maybe it finds non-linear effects: “If student spent < 2 minutes on the article, virtually zero chance of learning; if > 5 minutes, then decent.” That could lead to setting a minimum time threshold for awarding any proficiency boost. - We can also use ML to identify at-risk students early: e.g., cluster engagement patterns to see who is likely to fail despite content usage, and then intervene (this is more analytics than proficiency measure, but related).

- Eventually, the ML model could even replace some manual rules: for instance, instead of computing engagementScore via a fixed formula, we could train a model that inputs raw metrics and outputs a “predicted learning gain” value. If we trust it and it’s validated, that could be more accurate. But it should still work within the bounds of the Bayesian framework (i.e., output an increment that doesn’t exceed logical limits).

7. Decay (Forgetting) Model: - Incorporate a forgetting mechanism: If a student learned something only via content and then weeks pass with no practice, gradually decrease the proficiency estimate. Duolingo’s half-life regression can guide how to do this ²⁷ ²⁸ – basically, reduce probability as a function of time, faster if no reinforcement. - This can be skill-specific: each skill could have a half-life based on difficulty or the student’s profile (for example, math procedures might be retained for X weeks, pure memorization facts might decay faster without rehearsal). - The system can counteract decay by suggesting review content or quick quizzes. If the student re-engages (reads again or practices), the probability can be boosted back.

8. Reporting and Interface: - Present proficiency in a user-friendly way. Possibly a **dashboard with skill bars** (like Duolingo’s strength bars, or Khan’s mastery levels but updated with content learning). For

example, a bar might say "Algebraic Equations: 60% mastery" and if that 60% is mostly from content study, it might have an icon or color indicating "the system infers you might know this, but you haven't proved it yet." - Provide an **explanation feature** (especially for instructors): On clicking a skill, they can see "Mastery sources: Section 3.1 reading (+20%), Video: Solving Equations (+15%), Practice Quiz 1 (+30%, confirmed)." This transparency builds trust and helps instructors identify if a student has been learning passively or actively. - If a student's proficiency is lagging, the system might highlight which content they spent little time on or skipped, as actionable feedback ("Consider reviewing topic X's reading, since you haven't spent much time on it and your practice scores are low").

9. Validation Plan: - Before full rollout, pilot the framework on historical data if available. For example, take a past course where students had reading assignments and quizzes. Simulate the proficiency estimation for those students (would require reconstructing logs) and see how well it would have predicted quiz outcomes or final exam scores. Adjust parameters accordingly. - Alternatively, run an A/B test in a live class: one group uses the new proficiency estimator visible (and maybe gets recommendations based on it), another group doesn't. Check if the group with it performs better or is more engaged, etc. Also check if the proficiency estimates align with actual graded performance at course end. - Continuously gather user feedback: do students and teachers find the proficiency reports accurate? This qualitative feedback can catch situations like "it says I mastered topic X but on the midterm I did poorly on it, so something's off."

Implementation Complexity and Feasibility: - Much of the needed data (time, clicks) is already captured by online learning environments. The main implementation work is developing the logic to compute and update the proficiency states. This can be done as a service within Nexus AI that subscribes to event streams. - The Bayesian update formulas are not computationally heavy; they can run in real-time for each event. The data storage (per-student skill probabilities) is manageable. - The trickier parts are content-skill mapping (needs domain expert input and maybe alignment with standards) and initial parameter setting. We might use defaults from literature: e.g., reading a textbook section might on average have a 0.3 probability of teaching a new concept (30% learn rate), watching a video might be similar (~0.25–0.4 depending on interactivity), doing a practice question if answered correctly is nearly proof (we'd set that effectively as confirming mastery with maybe 0.8+ probability). - Over time, the system should override these with learned values from actual usage data as described.

Example Scenario (Illustrative): Consider a high school student, Jane, learning quadratic equations: - She reads the textbook section on "Derivation of quadratic formula" for 10 minutes (expected 8 minutes, she also scrolls up and down reviewing some steps). engagementScore comes out to 0.9 (she thoroughly read it). That section is mapped to skill "Quadratic formula derivation" with $L_c = 0.4$. Jane's prior on that skill was 0.1 (she's brand new to it). Update: new $P \sim 0.1 + (0.9 * 0.4 * 0.9) = \sim 0.46$ (46%). Significant jump, but she's not near mastery yet. - Next, she watches a Khan Academy video on the same topic on YouTube (suppose Nexus notices this via some integration or she marked it as supplemental resource). The video is 5 minutes; she watched fully and replayed one part – engagementScore 1.0. Video has $L_c = 0.3$ for that skill (reinforcement). Now $P \sim 0.46 + (1 - 0.46) * 0.3 * 1.0 \approx 0.67$ (67%). - Now she attempts a practice problem: "derive the quadratic formula from completing the square." She struggles and gets it wrong. The system has a slip model but basically, a wrong answer might drop her mastery by a certain amount (depending on model, maybe down to 40%). It now knows content alone wasn't enough; perhaps the system says: we estimated 67%, but the error indicates a gap. It might now explicitly recommend: "Review the derivation again or see another explanation, then try a similar problem." - She goes back and reads a worked solution example provided (content item with $L_c = 0.2$, engagement 1.0 since she reads it fully). P increases: $0.4 + (1 - 0.4) * 0.2 = 0.52$ (52%). Then she tries another problem and gets it right; now

mastery jumps to say 85%. - At this point, the skill is considered demonstrated. The earlier content-based increments helped get her there, but until that question was solved, the system maintained some uncertainty. Now it'll trust the 85% (maybe set to high confidence). - If she never got a chance to answer a question (say this was purely learned and the course moves on), the system would keep her at whatever probability and possibly decay it, but also mark it as unverified knowledge.

This scenario shows how content and questions interplay. Importantly, if Jane had aced the question on first try, the system might have spiked her to ~90% and there'd be no issue. If she kept failing despite content, the system would realize content wasn't sufficient and keep proficiency low until she succeeds or gets direct help.

Validation and Calibration

Ensuring the framework's accuracy and fairness is as important as building it. Key steps and strategies include:

- **Correlation with Outcomes:** We will regularly compute how well the estimated proficiencies correlate with independent measures:
 - Quiz and test scores (short-term).
 - Perhaps external exams or standardized tests if available (long-term).
 - Instructor evaluations or student self-assessments (qualitative checks).

Ideally, if our proficiency measure is good, a student with 80% mastery in our system should very likely score around 80% on a test of that material. If we find systematic bias (e.g., our 80% students are scoring 60%, then our model is overestimating and needs recalibration).

- **Item Response Theory (IRT) Calibration:** We can use IRT on assessment questions to fine-tune our proficiency scale. If our system predicts probability of correctness, we can compare that to actual probability observed. For instance, among all instances where our model said "student has 0.7 probability on skill X", did about 70% actually get it right? If not, adjust.
- **A/B Testing Interventions:** The ultimate test of proficiency estimation is whether acting on it improves learning. We can try interventions like:
 - For students whose content-based mastery is high but unconfirmed, send half of them a prompt to do a quick practice quiz and see if it secures their learning vs. the half who aren't prompted. If the prompted group ends up doing better, it validates that we successfully identified a need and addressed it.
 - Use the estimates to personalize study plans (e.g., a review module for those with low proficiency), and measure if that closes gaps.
- **User Studies:** Interview or survey instructors using the system: does the proficiency info align with their perception of students? Do they spot any students where the system was very wrong? Collecting these anecdotes can lead to discovering corner cases (e.g., the system might be tricked by a student who left videos playing overnight – a teacher might flag "this student clearly doesn't know

the material despite the system showing mastery"). Those insights would lead to adding safeguards (like capping daily credit to reasonable limits, etc.).

- **Fairness and Bias Check:** We should ensure the model works equally well for different subgroups (age, language, etc.). For example, if English language learners spend more time reading (due to translation) but perform the same as native speakers, we might inadvertently give them more proficiency credit (for longer time) which might not translate to better test performance. We need to adjust so that time from slower reading due to language is not over-valued. This might involve normalizing engagement metrics per student (i.e., use each student as their own baseline – if this student usually needs 15 minutes to understand a page, then 15 is “normal” not “high”). Continual bias audits should be done.
- **Error Analysis:** Examine cases where the model was wrong (predicted high proficiency but student failed, or predicted low but student succeeded unexpectedly). Determine if there were telltale signs we missed:
 - Did the student perhaps get help outside (so they actually learned from another source that we didn’t track)? If so, that’s outside our system’s vision – an inherent limitation.
 - Did the student guess on the assessment? Maybe our model was right about proficiency but the question outcome was noise (this is why multiple data points are good).
 - Or did we mis-evaluate their engagement? For instance, maybe the student opened the content but didn’t really engage but our metrics counted it. We can tighten those heuristics.
- **Quality of Sources:** Use the data to evaluate content items. If some content consistently yields less learning than expected, maybe it needs improvement or replacement. Conversely, content with high learning efficacy could be highlighted or used more.

By iteratively calibrating, the framework will become more robust. We expect that in early phases, we will be conservative (not assigning overly high mastery from content alone) until we have evidence to increase those values. Over time, the confidence in content-based estimation should grow as we validate it.

Conclusion and Recommendations

In this deep research, we explored how proficiency gain can be estimated from non-question interactions and proposed a comprehensive framework synthesizing the best of heuristic, analytic, and machine learning approaches. To conclude:

- **Feasibility:** It is indeed feasible to glean significant insight into student learning from content engagement data. While it cannot fully replace traditional assessments, it can greatly enhance **continuous monitoring** and personalized support. Platforms like Khan Academy and Duolingo implicitly acknowledge this by encouraging engagement and correlating it with outcomes, even if they don’t formally award mastery for it. Pearson’s Nexus AI can take a pioneering step by quantifying these insights in a transparent, data-driven way.
- **Implementation Priority:** Start with **infrastructure** – ensure content is well-tagged and interaction data is captured in detail. Then implement the **Bayesian update model with heuristic parameters**

(this will give an immediate working system). Simultaneously, design the **user interface** for showing proficiency so that it's intuitive (perhaps using terms like "learning progress" or "estimated mastery" with tooltips explaining it's based on activities, not just tests).

- **Short Term (Next 6-12 months):** Pilot the system in a controlled environment. Use an existing course with willing instructors to test it. Focus on refining the engagement scoring and L_c values. Ensure the system's suggestions (like prompting practice) are reasonable and not too frequent or too sparse.
- **Long Term:** Invest in **data analysis and ML** as more data comes in. Possibly form a small "learning analytics lab" team that continuously evaluates the efficacy of the proficiency estimates and works on model improvements (similar to how Duolingo has a research team for learning science). Keep humans in the loop for validation especially when deploying major changes.
- **Impact:** When fully realized, this framework will allow Pearson to deliver **real-time, granular insights** into student learning. Students will benefit from timely feedback ("I should revisit that section, my understanding isn't solid"), and instructors can perform targeted interventions ("these 5 students read everything but still don't get it—maybe they need a different approach, or they skipped the practice"). At a higher level, content creators can see how effective their materials are in teaching (beyond just people liking them, do they lead to mastery?). All of this aligns with Pearson's mission to improve learning outcomes with evidence-based, AI-driven tools.

In summary, our research strongly supports integrating non-question interaction data into proficiency evaluation. With careful modeling and validation, Pearson Nexus AI can lead the way in transforming passive learning data into actionable intelligence about student knowledge. This will make learning more personalized and assessment more continuous and seamless, ultimately driving better learning gains across diverse subjects and learners.

Annotated Bibliography

- **Gonzalez-Brenes et al. (2016) – "A Data-Driven Approach for Inferring Student Proficiency from Game Activity Logs":** (Peer-reviewed conference paper, ACM Learning at Scale 2016)
Quality: High. This paper is co-authored by Pearson researchers and academics, indicating rigor.
Summary: Introduces the SPRING framework for **stealth assessment** using educational game log data. The authors eliminated the need for expert-defined models by using machine learning to interpret sequences of in-game actions. They achieved a correlation of 0.55 with traditional test scores ⁶, validating that non-question behaviors can predict performance. This source underpins our section on using ML and log data for proficiency inference, demonstrating feasibility in a Pearson context. It also provides insight into Evidence-Centered Design and the costs of engineering evidence models ¹⁰. We used it to exemplify data-driven modeling and to support the idea that rich log data can partially replace explicit tests.
- **Yürüm et al. (2022) – "The use of video clickstream data to predict test performance":** (Peer-reviewed journal article, Education and Information Technologies, Springer)
Quality: High. It's a recent study with solid methodology (two experiments, data mining techniques).
Summary: Examines how fine-grained video interaction behaviors (pauses, rewinds, fast-forwards)

correlate with and predict university students' test scores. Found that *more interactive engagement (pausing, rewinding) correlates with higher performance*, and they could predict scores with RMSE 15-20% ². Important features were backward seek speed and number of pauses ². We cited this to emphasize the importance of **engagement quality** and to justify our focus on those metrics (not just time). It reinforces that passive vs. active video watching have different learning outcomes, guiding our engagement scoring system.

- **Okayama et al. (2024) – "Identifying Key Indicators of Proficiency in Junior High Math: Roles of Daily Handwriting Learning Logs":** (Conference paper, ICCE 2024 – likely peer-reviewed)

Quality: Medium-High. A niche study, but within a reputable conference.

Summary: Investigates whether **handwritten notes and work** (captured via a digital system) contain signals of proficiency. Found weak but significant correlations between features like writing speed, task time, and proficiency ¹². It suggests even subtle behaviors have some relationship with learning. We used this to support the idea that *learning process data (like note-taking behavior)* can inform proficiency, albeit modestly. It adds depth to our argument that multi-modal data (not just clicking but writing/drawing interactions) can be harnessed in our framework when available.

- **Khan Academy Research Brief (2023) – "Every minute on Khan Academy leads to learning gains":** (Efficacy report/blog by Khan Academy research team)

Quality: Medium. It's not peer-reviewed, but it's data-driven and authored by Khan's in-house researchers, referencing a study with 100k students.

Summary: Used a within-student analysis to show increased Khan Academy usage from one year to the next yielded improved MAP Growth test gains ¹. They used a Conditional Growth Index to measure this. We referenced this to establish that *engagement time has measurable effects on learning outcomes*, lending credibility to our assumption that more content interaction can equate to proficiency gains. It also provided a real-world efficacy proof and an example of correlational analysis in a broad setting.

- **Khan Academy Support Thread (2020) – "Videos not mandatory for mastery":** (Khan Academy Help Center forum, answer by staff)

Quality: Low (informal source), but authoritative for Khan's features.

Summary: Confirms that Khan's mastery system only counts exercises/quizzes, *videos do not count toward mastery* ¹⁵. We used this to illustrate how a major platform treats content vs assessment, highlighting the conservative approach in current systems. It was useful evidence in case studies, emphasizing the gap our framework intends to fill (by giving some credit to content).

- **Duolingo Blog (2019) – "How we learn how you learn":** (Official Duolingo blog post by staff, describing their Half-Life Regression model)

Quality: Medium. While a blog, it references a published paper (Settles & Meeder 2016) and explains the model in detail, indicating credible technical content.

Summary: Explains the development of Duolingo's HLR model for spaced repetition, combining forgetting curve theory with logistic regression on large data ¹⁶ ²⁵. It shows the formula and notes that HLR outperformed older methods by halving error rates ²⁶. We used this to discuss advanced modeling of proficiency (memory decay, practice effects) and as an example of ML improving a theoretical model. It supported our recommendations on implementing forgetting and using big data to tune learning models. It also provided an example of A/B testing an algorithm

change and measuring user retention improvement, which we found useful when suggesting validation techniques ¹⁸.

- **Tong & Ren (2025) – "Deep knowledge tracing and cognitive load estimation for personalized learning path":** (Peer-reviewed journal, Scientific Reports, Nature)

Quality: High. Scientific Reports is a reputable journal.

Summary: Proposes a dual-model combining Deep Knowledge Tracing (DKT) with cognitive load estimation (CLE) to optimize learning paths ¹³ ²¹. Reports improved prediction accuracy and learning efficiency in experiments. We didn't delve deeply into the technicalities in our report due to complexity, but we cited it to highlight the cutting edge: integrating *behavioral/cognitive state data* with knowledge models ¹⁴. It gave credence to our mention of cognitive load and the idea that future models could incorporate things like frustration or mental effort for better proficiency estimation.

- **Shute et al. (2011) – "Stealth Assessment in Computer-Based Games":** (Book chapter/paper by Valerie Shute, notable in this field)

Quality: High. Valerie Shute is a leader in stealth assessment, and her work is well-cited.

Summary: Although we didn't cite this directly above (it was in our search results), it generally covers using evidence-centered design to embed assessments into games unobtrusively. It likely provided background context for our understanding of stealth assessment and ECD. We indirectly referred to Shute's approach when discussing Bayesian networks and evidence models needing expert input ¹⁰. The concept of stealth assessment influenced our framework's philosophy: gather evidence invisibly during content interaction.

- **Early MOOC Prediction Study (2020) – in LNCS or similar via PMC**

Quality: Medium. Many MOOC studies exist; the one we accessed via PMC was open access and had a decent sample.

Summary: Focused on video interaction metrics (attendance, utilization, watch index) to predict who passes the MOOC ⁴. Showed that after just one week these metrics were effective predictors. We cited this to support early intervention and that *engagement data has predictive power very soon in a course*. It reinforces using content engagement as a diagnostic tool. Though we didn't deeply describe the specific metrics (they were formula given as images), we gleaned that essentially "how much and how thoroughly videos are watched early on" correlates to success. This backs our idea of using initial content interactions to adjust proficiency estimates and not waiting until after a failure occurs.

Each source above contributed to our understanding and justification of various components in the framework. The academic papers provided evidence of correlation and model efficacy, while industry examples gave practical viewpoints on implementation. We prioritized peer-reviewed research for claims about what works, and used industry documentation to ensure our recommendations align with what's viable and observed in real educational settings. All sources were checked for recency (most are within the last 5-10 years, ensuring up-to-date insights, with the exception of foundational concepts like stealth assessment which remain relevant).

- 1 New study finds every minute spent on Khan Academy can lead to learning gains - Khan Academy Blog
<https://blog.khanacademy.org/new-study-finds-every-minute-spent-on-khan-academy-can-lead-to-learning-gains/>
- 2 3 22 The use of video clickstream data to predict university students' test performance: A comprehensive educational data mining approach - PMC
<https://PMC9617048/>
- 4 5 Early Prediction of Success in MOOC from Video Interaction Features - PMC
<https://PMC7334715/>
- 6 7 8 9 10 11 20 23 24 A Data-Driven Approach for Inferring Student Proficiency from Game Activity Logs
https://people.cs.pitt.edu/~falakmasir/docs/gain_paper2016.pdf
- 12 library.apsce.net
<https://library.apsce.net/index.php/ICCE/article/download/5052/4989/6518>
- 13 14 21 Deep knowledge tracing and cognitive load estimation for personalized learning path generation using neural network architecture | Scientific Reports
https://www.nature.com/articles/s41598-025-10497-x?error=cookies_not_supported&code=5ad4b6ca-44db-4ae3-92f1-bef0153cff80
- 15 videos counting towards mastery & assignments towards mastery - Khan Academy Help Center
<https://support.khanacademy.org/hc/en-us/community/posts/360073652392-videos-counting-towards-mastery-assignments-towards-mastery>
- 16 17 18 25 26 27 28 Duolingo Blog
<https://blog.duolingo.com/how-we-learn-how-you-learn/>
- 19 [PDF] Learn Anywhere - English Performance Standards - Pearson
<https://www.pearson.com/content/dam/global-store/global/resources/efficacy/evidence-about-learning/Pearson-Learning-Design-Principles-English-Performance-Standards-summary.pdf>