# Adaptive Assessment Psychometrics Across Developmental Stages

Behavioral inference-based adaptive assessment spanning ages 5 through adulthood is technically feasible but requires fundamentally different psychometric approaches at each developmental stage. The most critical implementation decision involves **model selection by age**: Rasch/1PL models for early childhood (5-7), transitioning to 2PL for school-age and beyond, with multi-dimensional IRT supporting percentage-based multi-style profiles. However, a significant finding challenges the core premise: **the learning styles matching hypothesis lacks empirical support**, meaning behavioral detection of preferences is measurable and reliable, but adaptation based on detected "styles" does not improve learning outcomes. This report provides actionable guidance for implementing developmental adaptive assessment while navigating this constraint.

---

## IRT parameters require age-specific calibration strategies

Item Response Theory calibration presents distinct challenges across developmental stages, with the most substantial differences occurring between pre-literate children and older populations.

**Discrimination parameters (a)** show greater instability in young children due to developmental variability in response patterns. NWEA MAP Growth addresses this by using a Rasch (1PL) model across K-12, (NWEA) (Intensiveintervention) which assumes equal discrimination—a pragmatic choice that produces stable estimates across developmental levels. For ages 5-7, target discrimination values of **a = 0.8-2.5** are recommended; values below 0.5 indicate weak items unsuitable for adaptive use.

**Difficulty parameters (b)** must account for construct-irrelevant factors that inflate apparent difficulty for young children: literacy demands (reading requirements artificially increase difficulty for pre-readers), cognitive load from complex instructions, and attention limitations that make longer items harder regardless of content. The NWEA approach uses a **grade-independent RIT scale (100-350)** where difficulty parameters span multiple grades, with extensive audio support for K-2 to remove reading confounds. (NWEA)

**Guessing parameters (c)** present particular problems with young children. Research on Raven's Colored Progressive Matrices found significant guessing behavior in children ages 5-11, especially among younger participants. (PubMed Central) The recommendation for ages 5-7 is to **avoid 3PL models entirely** and instead use open-ended response formats, 1PL/2PL models with items designed to minimize guessing, or performance-based observational items.

| Age Group | Recommended Model | Minimum Sample Size | Key Constraint |
| --- | --- | --- | --- |
| 5-7 | Rasch/1PL or GRM | 500 | Audio support required, 3-point scales maximum |
| 8-12 | 2PL or GRM | 500-1,000 | 4-5 point Likert viable |
| 13-17 | 2PL or GRM | 500-1,000 | Standard administration |
| 18+ | 2PL or GRM | 500-1,000 | Full complexity |

## Differential item functioning identifies which items require stage-specific calibration

Research on PROMIS Parent Proxy Report Scales (ages 5-7 vs. 8-17) found that the **majority of items functioned similarly across age groups**, (PubMed Central) but specific item types exhibited significant DIF. Items referencing age-specific activities (homework, sports teams), abstract concepts (stress, anxiety), and social contexts that change with development consistently showed measurement non-invariance. (PubMed Central) Notably, **virtually no DIF** was observed between ages 8-12 and 13-17, (PubMed Central) suggesting these can often share calibration pools with appropriate linking.

Items likely to be **age-invariant** and suitable for cross-developmental anchoring include concrete behavioral descriptions, simple preference items, items about basic emotional states, pictorial response items, and items referencing universal experiences. The recommended approach is a **hybrid calibration strategy**: develop age-specific item pools with overlapping content ranges, include **30-40% of items as cross-age anchors** between adjacent stages and **20% as universal anchors** appropriate for all ages.

For DIF detection, the recommended methods are IRT-based likelihood ratio testing, Mantel-Haenszel procedures, or MIMIC models, (Columbia University Mailman …) with a **purification process** iterating detection to avoid false positives. Effect size interpretation follows established thresholds: negligible ($\Delta R^2 < 0.035$), moderate (0.035-0.07), and large ($\geq 0.07$).

## Vertical scaling enables growth tracking across developmental transitions

Large-scale assessments provide validated models for vertical scaling. NWEA MAP Growth's approach demonstrates key principles: use of a single **Rasch model for stability across grades K-12**, (Mapnwea) a grade-independent scale (RIT units from 100-350), (NWEA) common-item linking between adjacent grade bands, and audio support for early grades to ensure construct equivalence.

The **Non-Equivalent Anchor Test (NEAT) design** with common items is recommended for developmental assessment, requiring **minimum 20-25% common items** between adjacent stages. Anchor items must span the

difficulty range of both groups being linked, show no significant DIF across age groups, (PubMed Central) demonstrate moderate difficulty to avoid floor/ceiling effects, and have high discrimination ($a > 1.0$). (ERIC)

For linking methods, separate calibration with **Stocking-Lord or Haebara linking** is recommended when age groups differ substantially in ability, as is typical across developmental stages. (PubMed Central) Concurrent calibration can be used when model fit is good and group differences are moderate. Research shows **Haberman joint linking** minimizes bias and RMSE compared to chain linking approaches. (Taylor & Francis Online)

A critical validity consideration for learning style assessment: researchers must determine whether "learning style" represents the same construct across ages 5-18+. **Bifactor models** can separate a general factor from age-specific factors, helping address potential construct shift across development.

---

## Continuous profile estimation requires Bayesian sequential updating

Traditional single-administration psychometrics don't apply to continuously evolving profiles. The **Sequential Bayesian (SB) method** represents the current state-of-the-art for continuous ability estimation, computing posterior ability at time t as:

**$\theta_t|t = \theta_t|t\text{-}1 + K_t \times r_t|t\text{-}1$**

where $K_t$ is the Kalman gain matrix and $r_t|t\text{-}1$ is the prediction error. Research by Xiong, Cohen & Xiong (2023) found SB methods showed more accurate and reliable estimation than concurrent Bayesian methods, particularly with smaller samples (N=150-500).

The **Kalman filter framework** provides optimal state estimation for evolving traits, distinguishing measurement error from true state changes, handling missing data without imputation, and quantifying uncertainty through maintained covariance matrices. This approach has been validated for tracking psychological constructs over time, including mood, ability, and preference profiles.

**Reliability for dynamic profiles** cannot use traditional test-retest or internal consistency approaches. Instead, **conditional reliability coefficients** are required:

**$\rho(X,X'|\theta) = I(X,\theta) / [I(X,\theta) + 1]$**

where $I(X,\theta)$ is the score information function. (PubMed) Traditional conditional standard errors alone are insufficient—research by Nicewander (2018, 2019) demonstrated that portions of score distributions where scores are most/least precise can be misidentified without conditional reliability analysis.

For validation of behavioral inference from observed actions to latent style preferences, **multitrait-multimethod (MTMM) designs** with explicit theoretical frameworks are required. However, meta-analysis found only r = .23 correlation between teacher ratings and observations of similar behaviors, (Sage Journals) suggesting careful construct definition is essential before behavioral inference.

---

## Multi-dimensional CAT for style profiles differs fundamentally from ability assessment

Computerized Adaptive Testing for non-cognitive constructs faces challenges absent in ability assessment: no "correct" answers, potential need for ideal-point rather than monotonic item response functions, high vulnerability to faking, and significant response style confounds (acquiescence, social desirability, extreme responding).

**Multi-dimensional IRT (MIRT)** supports percentage-based multi-style profiles through between-item MIRT where each item measures one primary dimension with correlated dimensions allowed, or **bifactor MIRT** with a general factor plus orthogonal specific factors. Research on bifactor MIRT-based CAT achieved **r = 0.94 correlation with full item bank scores using only 11 items**.

For **item selection algorithms**, the optimal choice depends on measurement goals:

- **D-optimality** (maximizes determinant of Fisher information matrix): Best when all dimensions are intentional targets

- **A-optimality** (minimizes trace of inverse information): Slightly outperforms D-optimality when all abilities are intentional; more computationally efficient

- **Kullback-Leibler information**: More robust at early CAT stages when estimates are imprecise; outperforms Fisher-based methods initially

- **Mutual information**: Yields smallest conditional bias across ability levels in simulation studies

For style assessment specifically, **precision targets can be lower** than ability assessment given typically lower stakes. Standard error thresholds of **SE ≤ 0.45-0.55** (equivalent to reliability ~0.70-0.80) are often sufficient, compared to SE ≤ 0.30 (reliability ~0.91) for high-stakes ability decisions. The recommended approach: use KL information early-stage, switch to A-optimality or D-optimality after 5+ items, with marginal SE thresholds per dimension for stopping.

**Faking controls** are essential for non-cognitive CAT: forced-choice formats with Thurstonian IRT models recover normative scores from ipsative data, ideal-point IRT models (GGUM) better fit preference data, person-fit statistics detect aberrant responding, and mixture IRT models identify response-style subgroups.

---

## Cognitive constraints dictate radically different assessment structures by age

Assessment design must accommodate developmental realities that fundamentally differ across age groups.

**Working memory capacity** increases from approximately **1.5-3 visual items** in young children to 3-4 items approaching adult levels by ages 6-8. (Springer) Verbal working memory shows about 4 chunks in young children versus 6 in adults, (PubMed Central) with children shifting from visual to verbal processing reliance around age 9. (ERIC) This constrains item complexity and response option counts for younger children.

**Attention span norms** follow a rule of thumb of 2-5 minutes per year of age: 5-year-olds sustain 10-25 minutes, (Brain Balance Centers) while adolescents can maintain focus for 28-48 minutes. (Discovery ABA) Research confirms attention span is significantly longer in young adults than children. (PubMed Central) For assessment design:

| Age | Maximum Session | Item Complexity | Response Format |
| --- | --- | --- | --- |
| 5-7 | 10-15 minutes | Single-step, concrete | Pointing, manipulating, verbal with visual support |
| 8-12 | 20-30 minutes | Multi-step with scaffolding | Mixed selection and short constructed response |
| 13-17 | 45-60 minutes | Multi-dimensional, abstract | Extended responses, complex performance tasks |
| 18+ | 60+ minutes | Full complexity | All formats appropriate |

**Pre-literate assessment (ages 5-7)** requires oral administration of all instructions, familiar and comfortable settings, picture-based or play-based tasks, individual administration, brief sessions (10-15 minutes maximum), and large colorful illustrations. (Children's Hospital New Orleans) (SimplePractice) Established tools like the **Preschool Language Scales-5**, **Bracken Basic Concept Scale**, and storybook-embedded assessments like PELI provide validated approaches.

---

## Gamified assessment can achieve traditional psychometric standards

Evidence on gamified assessment validity is more positive than often assumed. Landers et al. found convergence between latent game performance and latent cognitive ability of **r = 0.97** in a theory-driven game-based assessment, providing "strong evidence that traditional game development practices can be applied to create an assessment reaching traditional psychometric standards."

Additional validation evidence includes: stealth assessments in Physics Playground showing valid correlations with external physics test scores; game-based persistence measures in Poptropica achieving **reliability alpha of .87** with children ages 6-14; and deep learning stealth assessment predicting posttest scores with **effect size d = 0.89**.

However, **construct-irrelevant variance remains a concern**. Game-based assessments may inadvertently measure cognitive ability due to high cognitive load, often also measure unintended personality traits, and can show variance related to gender that contaminates composite scores. Benefits of gamification include reduced test anxiety, reduced stereotype threat effects (with reduced performance gaps observed for women, non-white, and lower SES individuals), and enhanced engagement reducing dropout.

The **Evidence-Centered Design (ECD) framework** for stealth assessment specifies: a Competency Model defining targeted knowledge/skills/attributes, an Evidence Model delineating observable gameplay behaviors linked to competencies, and a Task Model describing characteristics that elicit evidence. (ed) Research

concludes that "evidence for digital game-based stealth assessment suggests that it is valid, reliable, and can be used to support learning." (ed) (ERIC)

---

## The learning styles hypothesis lacks empirical support—but preferences are measurable

A critical finding challenges the core premise of style-based adaptation: **comprehensive reviews find no evidence that matching instruction to detected learning preferences improves outcomes**. (Springer) Pashler et al. (2009) found "very few studies have even used an experimental methodology capable of testing the validity of learning styles applied to education," (Sage Journals) with meta-analyses showing the matching hypothesis effect size of $d = 0.04$ (essentially zero). (Springer)

The distinction is important: children **do have stable preferences** with established test-retest reliability. These preferences **do not predict** better outcomes when instruction matches preference. (PubMed Central) Learning style inventories show lack of construct validity with no correlation between different instruments measuring ostensibly the same styles. (Frontiers)

**Implementation implications**: Behavioral inference can reliably detect preference patterns, but adaptation decisions should focus on teaching adaptable learning strategies aligned with task complexity rather than matching to detected "styles." (Springer) The assessment system can provide valuable information about engagement patterns, metacognitive behaviors, and task approach preferences—but the adaptation logic should not assume that presenting content in a "preferred style" will improve learning.

**Metacognitive development** also constrains behavioral inference: metacognitive ability improves significantly with age during adolescence, is highest in late adolescence, and plateaus into adulthood. **Most growth occurs between ages 12 and 15**. For ages 5-7, observation-based methods are required rather than any reliance on self-monitoring; children at these ages can describe cognitive strategies and monitor uncertainty, but explicit metacognitive self-regulation is limited.

---

## Technical implementation requires specific tool selections and configurations

**Open-source IRT calibration tools** comparison for multi-stage developmental calibration:

**mirt (R)** is recommended as the primary tool, supporting 1PL through 4PL, graded response, multidimensional, and bifactor models. (CRAN) (ResearchGate) Key functions include (multipleGroup()) for multi-group calibration with invariance testing, (fixedCalib()) for fixed-item calibration during linking, and (fscores()) for EAP, MAP, MLE, and WLE ability estimation. EM algorithm is effective for 1-3 dimensions; MHRM (Metropolis-Hastings Robbins-Monro) is recommended for 3+ dimensions. (ResearchGate)

**TAM (R)** provides ConQuest-compatible syntax with built-in plausible value imputation and WLE scoring, making it suitable for large-scale PISA-style assessment applications. (RDocumentation) (Alexanderrobitzsch) **py-irt**

**(Python)** offers GPU acceleration via PyTorch for batch scoring millions of responses, though it currently lacks multidimensional and multi-group functionality. (arXiv) (INFORMS)

For **equating and linking**, the **equateIRT (R)** package provides Haebara and Stocking-Lord characteristic curve methods, chain equating for multiple forms, and standard error estimation for equating coefficients. (CRAN)

**Sample size requirements** are stage-dependent:

| Model | Minimum N | Optimal N | Developmental Consideration |
|---|---|---|---|
| Rasch/1PL | 200 | 500+ | Most stable for young children variability |
| 2PL | 500 | 1,000+ | Viable for ages 8+ |
| 3PL | 1,000 | 2,000+ | Not recommended for ages 5-7 |
| MIRT (2-3 dimensions) | 500 | 1,000+ | Cross-loadings require more data |

**Real-time ability estimation** should use Expected A Posteriori (EAP) with 21-41 quadrature points for production CAT (~1-2ms per update), switching to Weighted Likelihood Estimation (WLE) for final scoring to reduce bias. (Cogn-IQ) Duolingo reduced latency from 750ms to 14ms by rewriting their session generator in Scala; computational efficiency is essential for responsive adaptive systems.

## Industry implementations demonstrate proven technical architectures

**Duolingo English Test** provides a model for ML-integrated adaptive testing: ML/NLP-driven item difficulty estimation eliminates traditional pilot testing, Rasch model CAT operates across 25,000+ items in 11 difficulty bins, and soft scoring (probabilistic grades 0-1 instead of binary) enables nuanced CAT scoring. (duolingo) Psychometric results include internal consistency of 0.96, test-retest reliability of 0.80 at 30 days, and correlations with TOEFL (0.77) and IELTS (0.78). Item exposure rate of 0.10% far exceeds the 20% industry threshold. (duolingo)

**NWEA MAP Growth** demonstrates vertical scaling: Rasch model with RIT scale (100-350) spanning K-12, (NWEA) marginal reliability 0.93-0.96, and correlations with state tests of 0.75-0.85. (nwea) Common items link adjacent grade pools, with item parameter drift monitoring ensuring scale stability over time. (NWEA)

**Knewton's architecture** (now Wiley) used knowledge graphs with concept prerequisites, extended IRT with per-concept proficiency parameters, and time-varying proficiency incorporating learning/forgetting curves modeled as $R = e^{(-t/S)}$ where R is retention, S is memory strength, and t is time. (Knewton)

### Standards compliance requires comprehensive documentation

The **AERA/APA/NCME Standards for Educational and Psychological Testing (2014)** establish requirements for adaptive developmental assessment:

**Standard 5.15** requires documentation of linking accuracy and applicability limits when scores are linked across developmental stages. **Standard 5.16** mandates documentation of CAT procedures including item selection algorithms, ability estimation methods, and test termination criteria. **Standard 6.1** requires standardized administration procedures, with implications for documenting adaptive algorithm specifications.

For developmental assessment specifically, **Standard 12.4** requires documentation of the rationale for score interpretations across developmental levels, including justification of construct continuity, measurement invariance evidence, and ceiling/floor effect handling.

The 2025 Standards update (announced 2024) will address AI/ML in testing with enhanced coverage—relevant for systems using machine learning for item calibration or profile generation.

---

## Conclusion

Implementing adaptive learning style assessment across developmental stages is technically achievable with current psychometric methods and tools. The key technical decisions involve using **simpler IRT models (Rasch/1PL) for ages 5-7** with transition to 2PL for older populations, **Bayesian sequential updating** for continuous profile evolution, **multi-dimensional IRT with A-optimality or D-optimality item selection** for generating multi-style profiles, and **hybrid calibration strategies** with 20-30% anchor items linking developmental stages.

However, the evidence fundamentally challenges adaptation based on "learning styles": preferences can be reliably measured, but matching instruction to detected preferences does not improve learning outcomes. The assessment engine should focus on detecting engagement patterns, task approach behaviors, and metacognitive indicators rather than traditional learning style categories—and adaptation logic should leverage these signals for optimizing challenge level, pacing, and scaffolding rather than content presentation modality.

For ages 5-7, gamified stealth assessment with play-based tasks, audio instructions, and pictorial responses can achieve valid measurement while accommodating cognitive constraints. The system should observe behaviors rather than elicit self-report for metacognitive constructs with young children. Technical implementation using **mirt (R) for calibration** and **EAP estimation for real-time scoring** provides a validated pathway, with industry implementations from Duolingo and NWEA demonstrating production-ready architectures.