

Human-in-the-loop AI: Enhancing Assessment Item Generation in Medical Education

Xiaomei Song, PhD, Case Western Reserve University

Kate Weber, PhD, Case Western Reserve University

A. A list of Gen AI tools

ChatGPT	<ul style="list-style-type: none">• Generates MCQs, short answer, essay, matching, and case-based questions• Strong for clinical vignettes and reasoning-based questions• Highly customizable with prompts (e.g., OSCE, EOR, NBME-style items)• ChatGPT 4.0 fast for bulk question generation
Gemini	<ul style="list-style-type: none">• Structured question generation from documents (PDFs, notes, slides)• Strong for summarization → question conversion workflows• Integrates with Google ecosystem (Docs, Drive, Classroom)
Claude	<ul style="list-style-type: none">• High-quality long-form question writing• Strong at clinical reasoning and nuanced MCQ stems• Handles long documents very well (large context window)• Often preferred for assessment quality and clarity of wording
Copilot	<ul style="list-style-type: none">• Generates questions from Word, PowerPoint, and web content• Strong integration with Microsoft 365 (Teams, Forms, Word)• Useful for educators already in the Microsoft ecosystem
NotebookLM	<ul style="list-style-type: none">• Modifies uploaded materials into summaries, flashcards, and practice questions• Very strong for course-based exam preparation• Ideal for converting lecture notes into question banks
Custom GPTs / AI Builders (e.g., built on ChatGPT, Claude, Gemini, etc.)	<ul style="list-style-type: none">• Modifies uploaded materials into practice questions for formative purposes• Very strong for course-based question preparation• Ideal for converting lecture notes and materials into question banks

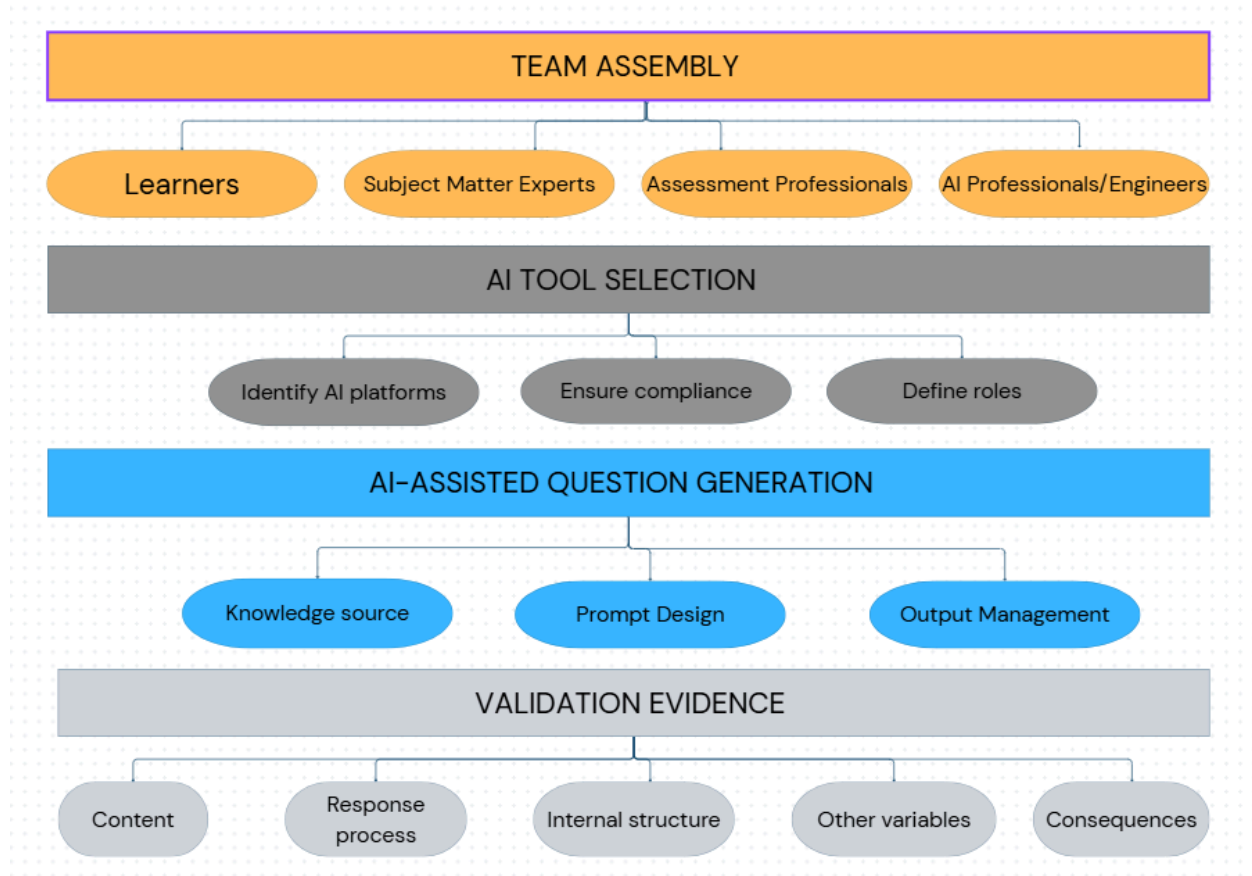
B. Draft Prompt

Text	Comment
<p>Context: You are an expert in evaluation methods preparing exam questions about Grant Writing.</p>	<p>Context helps the model frame its response in the terms you want. “You are a curious ten-year-old” would get completely different output</p>
<p>Using only the provided resources, generate one USMLE-style multiple-choice question for each learning objective in the 'Grant Writing LOs' document.</p>	<p>You are also telling it what information to use and giving it the overall task.</p>
<p>Question Structure: Clinical Vignette: Provide a realistic scenario. Lead-in: A clear, single question. Options: Exactly five (A-E). Key: The correct answer.</p>	<p>You explicitly define both the content and the format of the response.</p>
<p>Explanation: For each question, provide a 'Rationale' for the correct answer and a 'Distractor Analysis' explaining why each incorrect option is plausible but wrong. Identify the linked Learning Outcome.</p>	<p>The response must include justification from the model for its decisions to help the human reviewer assess the question.</p>
<p>Constraints: Focus on application and analysis (Bloom’s Taxonomy level 3+). Avoid 'all of the above' or 'none of the above.' Ensure the correct answer letter is randomized across the set.</p>	<p>Key constraints to shape how the questions are prepared</p>

Library:

- <https://gail.wharton.upenn.edu/prompt-library/>
- <https://mitsloanedtech.mit.edu/ai/basics/effective-prompts/>
- <https://structuredprompt.com/wp-content/uploads/2023/08/TRACI-Users-Guide-v1.pdf>

C. Human-in-the-loop AI Item Generation Flowchart and Checklist



Team Assembly	
<input type="checkbox"/>	Learners
<input type="checkbox"/>	Subject matter experts
<input type="checkbox"/>	Assessment professionals
<input type="checkbox"/>	AI professionals/engineers
Tool Design	
<input type="checkbox"/>	Identify AI platforms capable of question/item generation
<input type="checkbox"/>	Ensure data privacy, security, and institutional compliance
<input type="checkbox"/>	Define roles for AI vs. human review in workflow
AI-assisted Question Generation	
<input type="checkbox"/>	Knowledge source

<input type="checkbox"/>	Prompt design
<input type="checkbox"/>	Output management
Validation Evidence	
Content validity (see examples)	
<input type="checkbox"/>	Alignment: Ensure AI-assisted questions align with stated learning objectives and instructional content
<input type="checkbox"/>	Source Mapping/Representativeness: Confirm that each AI-assisted question clearly maps to a specific page or section in the provided material
<input type="checkbox"/>	Accuracy: Verify that each AI-assisted question is based on correct, evidence-based information (i.e., free from hallucinations or inaccuracies)
<input type="checkbox"/>	Difficulty Balance: Include a mix of both challenging and easy AI-assisted questions
<input type="checkbox"/>	Clinical Vignettes: Clinical scenarios may extend beyond the immediate topic; it is acceptable if questions can be answered without the vignette, as broader exposure supports learning
<input type="checkbox"/>	Distractor Quality: Ensure all AI-assisted question answer choices (distractors) are plausible and functionally relevant
<input type="checkbox"/>	Bias Review: Avoid using demographic characteristics (e.g., age, gender, race) as diagnostic shortcuts unless clinically justified
<input type="checkbox"/>	Images & Media: When an AI-assisted question references an image, ensure an appropriate, high-quality image is selected and inserted based on the description
<input type="checkbox"/>	Audit Trail: Document reviewer feedback, approvals, and edits for transparency and tracking
Response process (see examples)	
<input type="checkbox"/>	Think-Aloud Protocols: Ask students to verbalize their thought process while answering AI-assisted questions to identify misinterpretation or unintended cues
<input type="checkbox"/>	Cognitive Interviews: Conduct cognitive interviews with learners to confirm that their reasoning processes align with the intended difficulty and cognitive level
<input type="checkbox"/>	Reasoning Alignment: Confirm with students that correct responses require the targeted reasoning (e.g., application, analysis) rather than test-taking strategies or guesswork
Internal structure (see examples)	
<input type="checkbox"/>	Item Difficulty: Analyze item difficulty (e.g., p-values) to ensure an appropriate range across AI-assisted questions
<input type="checkbox"/>	Item Discrimination: Evaluate how well AI-assisted items differentiate between high- and low-performing learners (e.g., discrimination index, point-biserial correlation)
<input type="checkbox"/>	Reliability: Assess overall test reliability (e.g., Cronbach's alpha or KR-20 for dichotomous items)
<input type="checkbox"/>	Dimensionality: Evaluate whether AI-assisted items measure the intended construct(s) (e.g., factor analysis)
<input type="checkbox"/>	Score Consistency: Check for consistency of scores across years
Relationships with other variables (see examples)	
<input type="checkbox"/>	Criterion-Related Validity: Evaluate how well scores predict relevant outcomes (e.g., clinical performance, course success)

<input type="checkbox"/>	Convergent Evidence: Determine whether scores correlate with other measures of the same or similar constructs (e.g., course exams, standardized tests)
<input type="checkbox"/>	Discriminant Evidence: Confirm low correlation with unrelated constructs to ensure the test is not measuring unintended skills
<input type="checkbox"/>	External Benchmarks: Compare performance to established standards or norms when available
Consequences (see examples)	
<input type="checkbox"/>	Intended/unintended Impact: Evaluate whether the assessment supports or hinders learning, retention, teaching and desired educational outcomes
<input type="checkbox"/>	Fairness: Ensure the assessment does not advantage or disadvantage specific groups and supports equitable evaluation
<input type="checkbox"/>	Decision Accuracy: Review whether decisions based on scores (e.g., pass/fail, remediation) are appropriate and justified

Verbatim Suggested Prompt

You are an expert in assessment of medical students and development of multiple-choice exams. Using only the provided resources, generate one USMLE-style multiple-choice question for each learning objective in the 'Grant Writing LOs' document.

Structure:

Clinical Vignette: Provide a realistic scenario.

Stem: A clear, single question.

Options: Exactly five (A-E).

Key: The correct answer.

Explanation: For each question, provide a 'Rationale' for the correct answer and a 'Distractor Analysis' explaining why each incorrect option is plausible but wrong. Identify the linked Learning Outcome.

Constraints: Focus on application and analysis (Bloom's Taxonomy level 3+). Avoid 'all of the above' or 'none of the above.' Ensure the correct answer letter is randomized across the set.