

Feedback Generation in Education using Large Language Models: A Survey of Recent Advances

Anonymous ACL submission

Abstract

Large language models (LLMs) have made it possible to generate formative educational feedback at scale, but naïve generation often fails to meet educational requirements. A rapidly growing line of work reframes feedback generation as a decision-making problem: systems think *pedagogically* before they speak, by making intermediate choices—what to target, what action to take, how much support to provide, and what evidence should justify the message—before realizing the feedback. This survey reviews recent 58 papers on LLM-based educational feedback generation that incorporates such deliberative structure. We organize existing systems by where decision-making lives: prompting and in-context planning, training-time alignment, inference-time candidate selection, and scripted pedagogical scaffolds. We highlight open challenges in signal and judgement reliability, construct validity, generalization, personalization and scaling. We conclude with recommendations for building more auditable, controllable, and pedagogically-grounded feedback systems.

1 Introduction

Feedback is widely regarded as one of the most powerful influences on learning and achievement. Done well, it can substantially boost student performance; done poorly, it can hinder it or even harm it (Hattie and Timperley, 2007; Mandouit and Hattie, 2023). *Formative feedback* is the information communicated to learners with the intention of modifying their thinking or behaviour to improve learning, and is most effective when it is timely, specific, credible, and supportive (Shute, 2008).

In parallel, rapid advances in the applications of large language models (LLMs) to education have opened new horizons for intelligent tutoring, automated assessment and content generation, and made it possible to automatically generate feedback at scale. Recent studies explore LLM feedback in educational settings, both when it succeeds and

where it still fails. For example, LLMs have been widely studied for assessing student writing on criteria such as accuracy, specificity, and tone (Dai et al., 2024; Steiss et al., 2024; Escalante et al., 2023). Technical domains like computer programming have also benefited from LLM feedback identifying reasoning errors but it still sometimes struggle with subtle mistakes (Silva and Costa, 2025). Works using GPT models in assessment, formative (during instruction) and summative (evaluation), suggest the timely, personalised feedback is achievable, but with uneven quality and questionable reliability and impact on learning (Krumsvik, 2025; Jovic et al., 2025), and significant risks to academic integrity, the potential erosion of core skills, and the risk that AI feedback reinforces misconceptions or provides misleading information (Hasanein and Sobaih, 2023; García-López et al., 2025). Hallucination is another risk, underscoring the need for more robust mechanisms to verify and constrain model outputs in high-stakes domains including education (Huang et al., 2025; Sahoo et al., 2024).

Together, these strands point to a tension: educational feedback is an intervention that must be correct, pedagogically appropriate, and calibrated to the learner, which LLMs might not always get right through prompting, leading to generic, misaligned, or simply wrong feedback. This has led to a growing line of work that treats LLM-based feedback generation as a **deliberative** decision making process rather than a generation process. In this survey, we look at 58 papers (appeared in 2024/2025) that employs LLMs in feedback generation for education, addressing questions like what to optimize, where does the decision making live, what signals can guide the process, and where do experts fit in the loop.

1.1 Framing and Terminology

Feedback generation is an *instructional intervention* constrained by multiple objectives: it must be

content-correct, pedagogically aligned (e.g., hinting rather than revealing solutions), actionable, and appropriately calibrated in tone and confidence (e.g., not over- or under-praising). Feedback generation systems make at least one **intermediate decision** explicit before committing to the final output; typically selecting a feedback **target** (which error, misconception, or rubric dimension to address), an **action** (e.g., probe, hint, validate, request explanation, suggest revision, escalate), and **timing**. The feedback’s *grounding signal* is an observable, external cue that anchors feedback planning, including unit test results, revised answer quality or instructor edits (Ma et al., 2025; Wang et al., 2025a; Mok et al., 2025).

Thus, a feedback policy

$$\pi_F : (I, T, C, S) \rightarrow (A_F, t_F)$$

maps a learner’s input (I) for a task specification (T) under the learning context, curriculum and constraints (C) and the observable signals (S) into a timed feedback action (A_F, t_F). This framing allows comparing systems across domains by asking what learner artifact they utilize, what signals guide those decisions, and which constraints dominate.

1.2 Organisation

We organise this survey following a thematic flow: we discuss the domains and subject areas covered in the literature (§2), the architectures, mechanisms and workflows that are the backend of the surveyed systems (§3), then we discuss the gaps, open problems and opportunities that we observe from the pool of literature (§4). We also discuss evaluation strategies and rubrics in the Appendix (§C).

2 Domains and subject areas

2.1 Math

The works that handle math education in this survey (Table 1) concentrate in three settings: (i) *dialogue tutoring*, where planning is sequential and stateful, and feedback generation conditions on the full interaction to choose what is appropriate next with success predictors (Scarlatos et al., 2025b); (ii) *multi-choice questions setting*, where distractors provide explicit error handles, and feedback targets the misconceptions implied by, or explicitly associated with, the selected distractor (Scarlatos et al., 2024; McNichols et al., 2024a); and (iii) *open-ended response setting*, where free-form student solutions must be analysed for misconception

targets, using reference/generated solution, and the intervention strength level should be chosen accordingly (McNichols et al., 2024b; Baral et al., 2024; Tonga et al., 2025; Pal Chowdhury et al., 2024).

2.2 Programming

Coding papers in the survey (Table 1) fall into three settings: (i) *outcome-grounded executable tasks*, where feedback should help pass code checks (e.g., unit tests), in single-turn (Chae et al., 2024) or multi-turn (Wang et al., 2025b) with appropriate type (explanation, hint, question) and strength (Lohr et al., 2025; Heickal and Lan, 2024); (ii) *classroom-wide workflows*, where feedback is generated at the aggregate level but get reviewed, edited, and dispatched by the instructor (Tang et al., 2025a,b), constrained by explicit course policies (Scholz et al., 2025); and (iii) *interactive learning interventions* where the system supports the learning process itself (e.g., proposing debugging hypotheses (Ma et al., 2024) or reflection on alternatives (Naik et al., 2024)), with explicit focus on when to intervene (Ghoochani et al., 2025).

2.3 Writing

Writing papers in our survey (Table 1) cluster into three settings: (i) *revision-grounded feedback*, where the quality of a draft and predicted revision improvement shape feedback (Nair et al., 2024), sometimes with hierarchical criteria (Yuan et al., 2024; Chamoun et al., 2024; Wang et al., 2025c) or consistency check across repeated submissions (Bhojan and Xin, 2025); (ii) *interactive cognitive support*, where feedback targets different cognitive stages of the creative process (e.g., ideation vs. evaluation), balancing being helpful with maintaining writer agency (Göldi et al., 2024); and (iii) *persona conditioning*, where feedback is a function of a writer-defined perspective that acts as an explicit lens with a stance that aligns with the writer’s intent (Benharrak et al., 2024).

2.4 Other domains

We also identified literature that address spoken language assessment, multimodal language tutoring and STEM (e.g. science assignments), see Table 1.

3 Mechanisms

A *feedback mechanism* is the concrete approach in which a system performs decision making before committing to feedback. Mechanisms differ both in *where* decision making happens, *which*

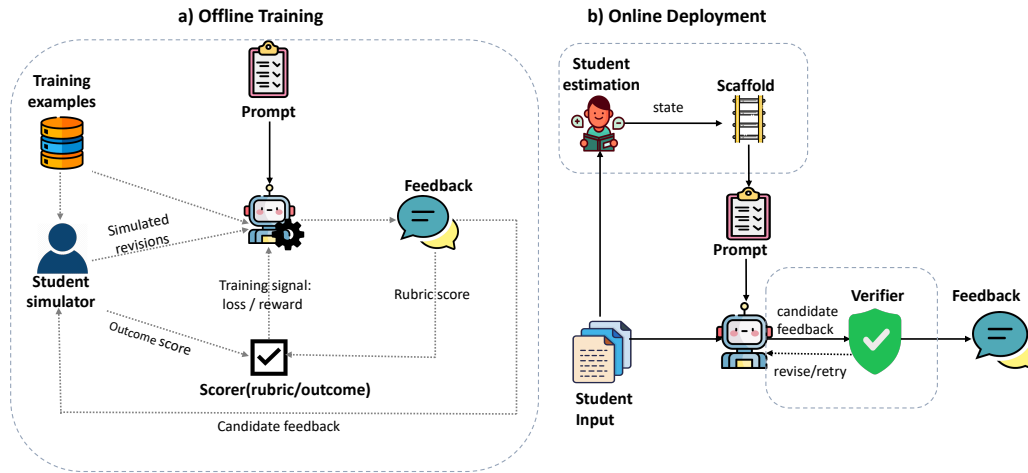


Figure 1: Overview of feedback generators: offline policy learning (Left), and scaffolded, verified feedback (Right).

intermediate decisions are available, *what* artifacts they introduce, and what supporting components they utilize. Note that these mechanisms are not mutually exclusive and can be synthesized.

3.1 Prompting and in-context learning

Here, the system uses an off-the-shelf LLM and describes the task and the constraints fully in the prompt. The feedback planning here happens within the model, selecting both (1) what to talk about (e.g., which error categories to address, or rubric) and (2) how (e.g., hint vs. explanation).

In writing, prompts explicitly guide the model to first elicit the relevant criteria and then generate critiques and suggestions that are traceable to them (Yuan et al., 2024), or to identify the most salient issues in drafts using a detailed rubric in zero- or few-shot settings (Rashkin et al., 2025), or using a rubric-based score (Han et al., 2024). In scientific writing, feedback pipelines similarly treat selecting what to focus on as a deliberate intermediate step, producing targeted revisions rather than broad, undirected commentary (Chamoun et al., 2024). Assessment workflow prompts produce multi-dimensional analytic judgments that can then be translated into actionable comments (Wang et al., 2025c). In math, in-context examples can help generate distractors and aligned explanations, constraining the feedback that follows (McNichols et al., 2024a). Finally, spoken-language assessment pipelines combine upstream transcription with rubric-aligned prompting to produce both scores and feedback (Shankar et al., 2025).

Prompting could also support personalization by injecting the right learner context (e.g., prior

attempts, misconceptions, or prior feedback) in the prompt. In (Reddig et al., 2025) the intermediate planning step is to choose the right context/examples to include to make the feedback align with the student’s needs.

3.2 Fine-tuning and preference learning

Here, an LLM is optimized at training time for feedback quality (*rubric-grounded*), or to the outcome after adopting the feedback (*outcome-grounded*).

Rubric-grounded preference optimization.

Several surveyed systems defined good feedback by an explicit pedagogical rubric, involving correctness, non-revelation of answers, diagnostic value, actionability and tone. Fine-tuning or preference learning is used at training-time to optimize a feedback policy toward ‘good’ feedback.

This pattern is explicit in diagnostic math feedback, where contrastive pairs of rubric-labeled feedback candidates are used for offline preference learning with DPO (Scarlatos et al., 2024). A closely related approach appears in pedagogical-alignment training, where structured pedagogical rubric, designed to distinguish guidance from answer giving, is used to build a preference dataset (Sonkar et al., 2024). Training-time optimization also appears in multi-agent settings (Yang et al., 2025) where collaborative agents are optimized for group discussions with pedagogy-guided interaction modelling, where intervention timing and strategy are modelled explicitly.

Outcome-grounded optimization A stronger form of objective-based decision making targets task success after adopting the feedback instead

of rubrics, usually requiring student simulation or explicit tool usage for verification at training time. For tasks with executable outcome, the grounding signal for feedback is the execution results after the simulated learner acts on the proposed feedback. (Chae et al., 2024) trains a feedback policy for the coding domain using Reinforcement Learning (RL) with an editor model (acting like a simulated student) that revises code given feedback, and unit test performance deciding rewards. In math tutoring, preference optimization is similarly used to push a tutor toward a hybrid objective that combines rubric-guided pedagogically-appropriate feedback with a prediction of learner success using knowledge tracing (Scarlatos et al., 2025b).

Beyond executable outcomes, (Dinucu-Jianu et al., 2025) frames the core tension of problem-solving vs. *teaching* problem-solving as a learnable policy, using RL to push tutoring toward pedagogically-desirable trajectories (e.g., reasoning). Synthetic data generation could also be utilized (Qian et al., 2025b), e.g. with simulated multi-agent teacher–student assignment–feedback loops (Zhang et al., 2025a), or multi-step synthetic Socratic dialogues between student, teacher and dean agents are used to fine-tune models for better pedagogical alignment (Liu et al., 2024a).

3.3 Runtime selection with learned verifiers

Another recurring mechanism is explicit inference-time feedback ranking of generated candidate feedback actions, evaluated by a learned scorer module. (Wang et al., 2025b) instantiates this mechanism in coding tutoring by sampling many candidate utterances for each turn, based on an estimated student state, then ranking them by a turn-by-turn scorer for the one most likely to make progress toward task completion. This is particularly useful when the action space is large (many plausible tutor moves).

LLM evaluators are also used to rank feedback candidates along pedagogical criteria in multi-stage agentic pipelines, e.g. in writing (Chamoun et al., 2024), or as a comprehensive automated evaluation layer that serves as a gatekeeper (Qian et al., 2025a). The reliability of these judges is a concern (Seo et al., 2025), but iterative evaluation and revision can stabilize feedback and reduce variance in reflective writing (Bhojan and Xin, 2025).

3.4 Scripted pedagogy scaffolds

Some of the surveyed works embed LLMs within external pedagogical scripts/ intervention policies

that decide what kind of support to provide when. The scaffold provides explicit decision making for the learning stage and trigger conditions, intervention timing, type and stance, before relaying that to the LLM. It is argued that these scripts should be grounded in assessment evidence and decision modules (e.g., ZPD-informed adaptation) to determine their effective timeline (Cohn et al., 2025).

For academic writing, (Chen et al., 2025) introduce an 11-step scaffold to guide the learner from identifying key questions to readability and grammar checking. Programming feedback ladders impose ordered levels of support (identification to progressively more directive hints) (Heickal and Lan, 2024; Lohr et al., 2025), as well as in math (Tonga et al., 2025). Alternatively, finite state transducers with solution space decomposition guides the process (Pal Chowdhury et al., 2024).

In the collaborative SQL optimisation domain (Naik et al., 2024), types of reflection get triggered by pattern matching, prompting ChatGPT to produce timely, personalised interventions that help surface alternative solution paths. For peer review writing, inspired by the Cognitive Process Theory of Writing, assistance is explicitly conditioned on cognitive goals (ideation vs. evaluation) and constrains support through interface design (Göldi et al., 2024). Persona-conditioned feedback makes *stance* a planned variable, letting writers define audience lenses that shape feedback perspective/specificity (Benharrak et al., 2024).

Alternative workflows that benefit from LLMs include *Learning-by-teaching* (Ma et al., 2024) where novices teach simulated students, with structured tasks that keep attention on constructing hypotheses and test suites. The *productive failure* approach (Puech et al., 2025) keeps students in a beneficial struggle regime before providing resolution. More generally, multimodal language tutoring can embed a script of “pedagogical instruction” to guide the tutoring and modality use (Liu et al., 2024b), and agentic tutors can hard-code Socratic guided questioning as a scaffold while grounding the interaction via retrieval (Knievel et al., 2025).

4 Open problems and opportunities

4.1 Signal reliability and judge validity

Educational feedback pipelines make intermediate choices using signals from rubrics, LLM-as-a-judge scores and mistake detectors. These signals can be noisy, biased, or brittle in ways that

346 materially affect downstream decisions, especially
347 in subjective dimensions like tone, positivity, and
348 perceived helpfulness. This could surface in rubric-
349 grounded optimization (Scarlatos et al., 2024), as
350 well as for automatic judges where claims about
351 improved tutor behaviour are hindered by the eval-
352 uators diverging from human judgments (Scarlatos
353 et al., 2025b), by bias and imperfect constraint
354 checking (Ferraz et al., 2024), or sensitivity to sur-
355 face forms such as verbosity (Seo et al., 2025).

356 A partial remedy is to shift to grounded veri-
357 fiable signals when possible as judging can miss
358 subtle issues (McNichols et al., 2024b; Baral et al.,
359 2024). In programming, tool-grounded rewards
360 (e.g. unit tests) reduce reliance on judges (Chae
361 et al., 2024). Tool-interactive critique offers a
362 primitive for grounding revisions in external evi-
363 dence (Gou et al., 2024), while constraint decom-
364 position turns rubrics into explicit constraints for
365 selective refinement (Ferraz et al., 2024).

366 A more general requirement is to treat the eval-
367 uator/judge as a *first-class* component, not a con-
368 venience scorer, requiring **calibration and audit-**
369 **ing**. Probing LLM evaluators shows that agree-
370 ment and accuracy depend on evaluator choice
371 and setup. Additional evaluators can act as *con-*
372 *trol signals* for quality gating and regeneration,
373 e.g. the *Dean* framing which places a special-
374 ized feedback evaluator upstream of delivery to
375 reject low-quality feedback and trigger regenera-
376 tion (Qian et al., 2025a). Without human-grounded
377 validation, adversarial/counterfactual stress tests,
378 or uncertainty reporting for borderline decisions,
379 pipelines risk *correlated failure modes*: generator
380 and evaluator sharing blind spots, and high scores
381 becoming self-referential rather than evidence of
382 pedagogical quality, as could be seen in analytic
383 writing assessment (Wang et al., 2025c).

384 Finally, reliability issues can also arise *before*
385 planning starts, as in spoken-language assessment
386 where errors and subgroup performance differences
387 can propagate into scoring and feedback (Shankar
388 et al., 2025), motivating **uncertainty-aware**
389 **pipelines** that carry confidence from upstream per-
390 ception into downstream decisions, and decision
391 signals that are both **meaningful** and **auditable** in
392 domains without crisp executability like writing.

393 4.2 Construct validity

394 Even when evaluation signals are reliable, they
395 may measure imprecise or ill-posed educational
396 constructs, e.g. unit tests capturing functional cor-

397 rectness in place of pedagogical gains in coding
398 learning (Chae et al., 2024), superficial rubric-
399 based revision improvements in essay writing that
400 have a weak connection to learning outcome (Nair
401 et al., 2024), or hallucinated or incomplete gen-
402 erated solution trees that could hinder math tutor-
403 ing (Pal Chowdhury et al., 2024).

404 In the learning-science perspective, ensuring
405 what systems optimize is meaningful calls for
406 feedback design that is aligned with instructional
407 goals and learner outcomes (Stamper et al., 2024).
408 Process-centred interventions provide exemplars
409 where the target is the learning behaviour itself,
410 such as in hypothesis construction for debug-
411 ging (Ma et al., 2024), or situated reflection that
412 supports exploring alternatives (Naik et al., 2024).
413 Similarly, a classroom randomized control trial
414 found that GPT-4-driven interactive homework im-
415 proves engagement and learning outcomes in gram-
416 mar (Vanzo et al., 2025).

417 Still, what is lacking is a scalable way to connect
418 intermediate decision quality (e.g. target selection,
419 hint choice, escalation timing) to *real, validated*
420 *outcomes* to inform feedback planning, without re-
421 lying on expensive classroom trials. Explaining
422 *why* gains occur and *for whom* require aligning
423 what planners optimize with behaviours that plau-
424 sibly drive learning (e.g., revision choices, practice
425 quality, persistence), and reporting behavioural out-
426 comes (what learners change vs. ignore), retention,
427 and transfer. A concrete opportunity is to standard-
428 ize reporting that connects feedback attributes to
429 learner actions (edits made, strategy shifts).

430 4.3 Generalisation beyond simulators/proxies

431 Gains of some feedback generation research comes
432 under simulated learners (e.g. dialogue knowl-
433 edge tracing (Scarlatos et al., 2025a)), or revisers
434 (e.g. simulated writing revisers (Nair et al., 2024)).
435 Two risks are attached to that: (i) *fidelity mismatch*
436 when the simulated behaviour diverges from the
437 real learning, and (ii) *proxy overfitting* when a sys-
438 tem fails to generalise beyond the simulation.

439 Furthermore, while the use of synthetic corpora
440 (e.g. SCALEFeedback for CS assignments (Qian
441 et al., 2025b)) and agentic simulations (e.g. SEFL
442 self-feedback loops (Zhang et al., 2025a)) are ex-
443 panding to counteract scarce classroom data, a
444 transfer risk arises: improvements may reflect syn-
445 thetic error patterns rather than authentic learning.
446 Partial remedy for this includes acquiring **human**
447 **evidence** that interventions survive reality, even

at small scale (e.g., debugging support (Ma et al., 2024)), and more systematic effort to **elicit simulation parameters from real trajectories** to calibrate simulators and stress-test parameter sensitivity in prompts and learner profiles.

4.4 Rubric-compliant generation

Rubric compliance at generation still lacks robustness, especially in long contexts or multi-turn settings. Even strong models often violate at least one constraint in real multi-constraint instructions (Ferraz et al., 2024), leading to phenomena like overpraise or over-inference (Guo et al., 2024) or drift in scaffolded-help behaviour (Sonkar et al., 2024).

A partial remedy is to enforce explicit criteria through structured critique or checkable tests. This makes **criterion selection and prioritization** a key problem: what to surface now under limited attention, as models can focus on low-hanging issues and miss the most salient problems otherwise (Rashkin et al., 2025). Another remedy is to ground critique with tools when possible: tool-interactive critique reduces hallucinations by grounding feedback in external verification (Gou et al., 2024), and outcome-grounded approaches make constraints become signals with suitable checkers (Chae et al., 2024). Yet, producing the intended *kind* of feedback remains unreliable (Lohr et al., 2025), motivating structured scaffolds such as feedback ladders (Heickal and Lan, 2024). These make strategy explicit (e.g., productive failure intent selection (Puech et al., 2025)).

Overall, the field still lacks a widely accepted **pedagogical constraint schema** that is simultaneously decomposable (Ferraz et al., 2024), checkable when possible (Chae et al., 2024), and aligned with learning-science constructs (hinting vs. telling, autonomy support, desirable difficulty), especially for long-form responses (Rashkin et al., 2025).

4.5 Overhead

Feedback mechanisms increase cost and latency through candidate sampling, verifier calls and iterative refinement. Long tool-interactive correction loops are one driver of overhead (Gou et al., 2024), but selective refinement can reduce unnecessary rewriting by focusing on detected violations (Ferraz et al., 2024). In classroom workflows, responsiveness also depends on human factors: instructor-facing systems mitigate workload via structured review and attention-aware batching (Tang et al., 2025b,a). Some approaches offload cost to offline

training (Chae et al., 2024) or caching (Pal Chowdhury et al., 2024), but deployment feasibility still depends on online constraints. And with more planning, multi-agent loops, and evaluator gating, quality can improve but costs rise (Chamoun et al., 2024; Zhang et al., 2025a; Qian et al., 2025a).

A key opportunity is **adaptive budgeting**: put compute on verification/regeneration only when stakes or uncertainty are high. Spoken-language assessment makes this concrete as upstream uncertainty can be estimated to trigger conservative feedback, extra checks, or human review (Shankar et al., 2025). Also, characterizing the *cost-quality frontier* is needed: when additional candidates, refinement, verifier calls, or human review are worth it under strict classroom latency and constraints.

4.6 Safe scaling

Scaling feedback at the classroom level introduces systemic risk: errors can propagate widely if aggregation/reuse is faulty. Instructor-validation workflows mitigate this by structuring review and propagation, but they also expose the need for additional safeguards when issue detection is imperfect (Tang et al., 2025b,a). Tutor CoPilot supports live tutors with only suggestions (Wang et al., 2024b) but raises research questions about how AI reshape tutor behavior and student autonomy. Propagation safeguards must go beyond correctness to audit non-leakage, course policy compliance, subgroup fairness, and preservation of instructional intent.

Recent deployments highlight the stakes: automated project feedback can reach many students quickly (Ghoochani et al., 2025), and interactive homework tutoring can be deployed at classroom scale with measurable learning gains (Vanzo et al., 2025). The open problems are **monitoring and rollback**: detecting recurring misconceptions or policy-violations, preventing repeated errors across cases, and notifying affected learners. Safe scaling also requires governance of data and uncertainty. Privacy- and sensitivity-aware handling of learner data (audio, reflective writing) constrains what can be stored, retrieved, or reused (Shankar et al., 2025; Bhojan and Xin, 2025; Jovic et al., 2025), and uncertainty-aware control suggests that low-confidence cases may warrant different scaling policies (e.g., reduced propagation).

4.7 Adaptivity and Personalization

Many feedback generators remain generic, and when personalization is introduced (knowledge

tracing, writer-defined personas, attention-driven interfaces), it raises evaluation and safety questions. Learning-science framings argue personalization should be grounded in theory and evaluated beyond subjective preference (Stamper et al., 2024). A theory of adaptive scaffolding argues that pedagogical agents should modulate support based on learner signals while avoiding over- and under-scaffolding (Cohn et al., 2025), while in-context personalized tutor feedback relies on inferring accurate learner state and misconceptions (Reddig et al., 2025). Systems operationalize this via adaptive platforms that plan learning paths and per-step support (Chudziak and Kostka, 2025).

Personalization also extends beyond individuals: timely intervention in collaborative discussions requires modeling group progress and interaction dynamics (Yang et al., 2025). However, deployment evidence suggests motivational and context-aware tailoring remains uneven even when students prefer automated feedback (Ghoochani et al., 2025).

Two under-developed research directions are **calibration** (detecting when personalization is off: over-hinting, under-challenging, mismatched tone) and **equity** (auditing persona-conditioned policies for stereotyping or differential treatment).

4.8 Weak comparability

Comparability between feedback generation systems is hard because objectives and criteria are not made commensurate. And while the design space is converging, the field still needs standardized **objective bundles** that settles primary objective, constraints, acceptable proxies and required robustness checks, and enforcing explicit reporting of evaluator calibration and uncertainty when LLM judges are used. TutorGym and MathTutorBench provide building blocks by offering structured testbeds for tutoring behavior and pedagogical capability (Weitekamp et al., 2025; Macina et al., 2025), while efforts like SCALEFeedback (Qian et al., 2025b) support reproducible training and evaluation.

5 Discussion and Conclusion

This survey reviewed recent work (2024–2025) of LLM-based educational feedback generation. A consistent pattern emerges: making *pedagogically valid and reliable choices*—what to address, when to intervene, how much to reveal, what evidence to justify—under real deployment constraints. This survey contributes a cross-domain analysis that

organizes mechanisms by *where* decision-making lives (prompting, training-time optimization, verification and selection, scaffolding).

A core view in the literature is moving from prompting to *feedback policies*: explicit mapping that connects observable grounding signals to feedback strategies, whether learned (e.g., preference optimization) or scripted (e.g., scaffolding rules), in a way that is auditable, stable, and safe at scale. These intermediate decisions make systems more comparable across domains. The surveyed work shows that even lightweight explicit structure—such as selecting criteria before drafting (Yuan et al., 2024; Chamoun et al., 2024), laddered levels of help (Heickal and Lan, 2024), or intent/state transitions that operationalize productive failure (Puech et al., 2025)—can convert otherwise underspecified “helpfulness” into controllable behavior.

Naïve feedback generation fails in high-stakes educational settings due to LLM propensity to violate some of the domain’s multi-objective constraints (correctness, non-revelation, actionability, calibration, tone), especially in long contexts/multi-turn settings (Ferraz et al., 2024; Sonkar et al., 2024). Decision-making approaches make the constraints explicit, enabling mechanisms such as candidate selection with verifiers (Wang et al., 2025b), self-critique (Yuan et al., 2024), or rubric-grounded preference learning (Scarlatos et al., 2024).

Across mechanisms (§3), surveyed systems can be abstracted into a common pipeline: 1) **Grounding signals**: observable cues that anchor the intervention (unit tests, revisions, rubrics, instructor edits, interaction traces). 2) **Intermediate decisions**: discrete choices made *because* of those signals (what to target, what action to take, what strength, what evidence to cite, whether to escalate or verify). 3) **Feedback realization**: the surface language form, constrained by the action and policy (hint vs. explanation, Socratic prompt, tone).

This decomposition helps reconcile seemingly different approaches. Training-time approaches learn a policy that maps signals to decisions (§3.2); verifier-based approaches decide among candidates (§3.3); and scripted scaffolds implement explicit pedagogical control policies (§3.4). Instructor-in-the-loop systems further externalize decision-making into workflow primitives such as triage, verification, and reuse (Tang et al., 2025b), which are central to reliability at scale.

Crucially, this view is actionable as a *design and reporting template*. It forces clarity about what

649 signals are trusted, which intermediate decisions
650 are exposed, and what realization constraints are
651 enforced. This highlights when a system is implic-
652 itly relying on subjective judgments versus when it
653 can anchor decisions in checkable outcomes.

654 This also highlights a practical implication:
655 when reliable grounding signals exist (e.g., unit
656 tests), systems can rely less on subjective evalua-
657 tion and more on verifiable outcomes (Chae et al.,
658 2024). When signals are inherently subjective or
659 weakly grounded (e.g., writing quality), systems
660 must invest in robust rubrics, evaluator calibra-
661 tion, and evidence-grounded critiques to avoid un-
662 founded feedback (Chamoun et al., 2024; Seo et al.,
663 2025). This leads to key takeaways:

664 - **Explicit intermediate decisions** (targets, actions,
665 scaffold level, escalation) are a high-leverage route
666 to auditability and controllability (Yuan et al., 2024;
667 Heickal and Lan, 2024; Puech et al., 2025).

668 - **Signal validity is the bottleneck**: as systems be-
669 come more policy-like, failure increasingly comes
670 from mis-calibrated decision signals, especially
671 judge signals (Seo et al., 2025; Qian et al., 2025a).

672 - **Evaluation must be decomposed** into action qual-
673 ity, feedback quality, and impact; otherwise gains
674 are easily misattributed to surface fluency (Chae
675 et al., 2024; Sonkar et al., 2024).

676 - **Cost-aware reliability** for real classrooms: verifi-
677 cation/planning must be budgeted adaptively, not
678 uniformly (Ferraz et al., 2024; Tang et al., 2025b).

679 - **Learning-science constructs** should increasingly
680 become decision variables, not only prompt phras-
681 ing (Hattie and Timperley, 2007; Shute, 2008).

682 Decision-making and verification increase cost
683 and latency, creating a tension between reliabil-
684 ity and feasibility (§4.5). Many systems therefore
685 implicitly assume either generous compute (e.g.
686 multi-stage agentic pipelines) or settings where de-
687 lays are acceptable. Classroom deployments, how-
688 ever, often require real-time responsiveness and ro-
689 bust safeguards against widespread mistake propa-
690 gation (Tang et al., 2025b). This elevates two prac-
691 tical research directions: 1) **adaptive budgeting**:
692 allocate verification and regeneration selectively
693 based on uncertainty, stakes, and the availability of
694 grounding signals, rather than uniformly (Ferraz
695 et al., 2024). Concretely, triggers can be disagree-
696 ment between candidates, missing or low-quality
697 grounding signals, or high-stakes course policy
698 cases (e.g., solution leakage); and 2) **monitoring**
699 **and rollback**: as feedback is reused or broadcast at
700 scale, systems need mechanisms to detect system-

701 atic failure patterns (e.g., recurring misconceptions,
702 policy violations, subgroup disparities) and prevent
703 repeated harm (Tang et al., 2025b).

704 Deployment realism also includes governance
705 constraints beyond latency: privacy and sensitivity
706 of learner data and equity risks that arise when
707 evaluator bias or personalization policies behave
708 differently across subgroups (Shankar et al., 2025;
709 Nazaretsky et al., 2025). These affect safe scaling
710 because they restrict what can be stored, reused,
711 and propagated, and they shape which “signals”
712 are acceptable as decision inputs.

713 Finally, the surveyed literature reinforces a
714 learning-science point that is easy to lose in NLP-
715 centric framing: educational feedback is an *in-*
716 *structional intervention*, not merely error correc-
717 tion (Hattie and Timperley, 2007; Shute, 2008).
718 Several works explicitly model the tension between
719 solving and teaching problem solving (Dinucu-
720 Jianu et al., 2025), or structure support around pro-
721 ductive struggle and scaffolding (Puech et al., 2025;
722 Cohn et al., 2025). These directions suggest that
723 the next wave of planning-based feedback research
724 should more directly encode theoretically grounded
725 constructs—e.g., autonomy support, desirable dif-
726 ficulty, and scaffold fading—as decision variables,
727 and evaluate them with behavioral outcomes (what
728 learners do next), not only perceived helpfulness.

729 Operationally, this means taking pedagogical
730 principles as *choices* the system must make: e.g.,
731 selecting a Socratic prompt vs. a direct hint (au-
732 tonomy support), delaying solution revelation until
733 evidence of impasse (desirable difficulty), and re-
734 ducing hint specificity as competence improves
735 (scaffold fading). Framed this way, learning-
736 science alignment becomes an implementable
737 agenda rather than an after-the-fact description.

738 **Conclusion.** LLMs have made feedback genera-
739 tion widely accessible, but they emphasize the *pol-*
740 *icy* question. Across the recent work, we observe
741 meaningful convergence on planning-based formu-
742 lations that leverage grounding signals, intermedi-
743 ate decisions, and explicit workflows to improve
744 reliability, controllability, and scalability. We iden-
745 tify open issues of trustworthy signals, construct va-
746 lidity, generalization, and safe scaling under class-
747 room constraints. Addressing these requires closer
748 integration between NLP mechanisms, educational
749 theory and evaluation practice, alongside standard-
750 ized reporting that makes objectives, signals, and
751 intermediate decisions explicit.

6 Limitations

Temporal scope. This survey focuses on work published in 2024–2025, with the goal of capturing the most recent wave of LLM-based feedback systems that explicitly incorporate planning or decision-making. As a result, we do not attempt to comprehensively cover earlier foundational literature on intelligent tutoring systems, automated feedback, or pre-LLM approaches, and we may also miss very recent papers that appeared after our cutoff.

Search methodology and coverage. Our paper collection relied on a manual, keyword-driven search and screening process, rather than a fully exhaustive systematic review. We used queries centered on LLMs and feedback/planning, including terms such as LLM-generated feedback, formative feedback, automated feedback, LLM tutor/tutoring, hint generation, Socratic tutoring, rubric-based feedback, programming feedback, writing feedback, feedback policy, verification/verifier, and LLM-as-a-judge. Since terminology differs across NLP, learning sciences, and HCI/EdTech communities, and because some relevant work is framed as assessment, scaffolding, revision support, or instructional interventions rather than “feedback”, our coverage is necessarily incomplete and may omit pertinent studies that use different descriptors or appear in adjacent venues.

Space constraints. ACL page limits require prioritizing the core synthesis and taxonomy in the main paper. We therefore place extended material such as expanded tables, evaluation details, dataset summaries, and additional examples or elaborations into the appendix. This improves readability of the main narrative but may reduce the detail available in the main text for some subsets of surveyed work.

References

Sami Baral, Eamon Worden, Wen-Chiang Lim, Zhuang Luo, Christopher Santorelli, and Ashish Gurung. 2024. [Automated assessment in math education: A comparative analysis of llms for open-ended responses](#). *educational data mining*, pages 732–737.

Karim Benharrak, Tim Zindulka, Florian Lehmann, Hendrik Heuer, and Daniel Buschek. 2024. [Writer-defined ai personas for on-demand feedback generation](#). In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pages 1–18.

Anand Bhojan and Tan Li Xin. 2025. [Reflexai: Optimizing LLMs for consistent and constructive feedback](#)

[in reflective writing](#). In *Proceedings of the 17th International Conference on Computer Supported Education (CSEDU 2025) - Volume 2*, pages 387–394. SciTePress.

Hyungjoo Chae, Taeyoon Kwon, Seungjun Moon, Yongho Song, Dongjin Kang, Kai Tzu-iunn Ong, Beong-woo Kwak, Seonghyeon Bae, Seung-won Hwang, and Jinyoung Yeo. 2024. [Coffee-gym: An environment for evaluating and improving natural language feedback on erroneous code](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 22503–22524, Miami, Florida, USA. Association for Computational Linguistics.

Eric Chamoun, Michael Schlichtkrull, and Andreas Vlachos. 2024. [Automated focused feedback generation for scientific writing assistance](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9742–9763.

Fumian Chen, Sotheava Veng, Joshua Wilson, Xiaoming Li, and Hui Fang. 2025. [Coachgpt: A scaffolding-based academic writing assistant](#). In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 4051–4055.

Jarosław A Chudziak and Adam Kostka. 2025. [Ai-powered math tutoring: Platform for personalized and adaptive education](#). In *International Conference on Artificial Intelligence in Education*, pages 462–469. Springer.

Clayton Cohn, Surya Rayala, Namrata Srivastava, Joyce Horn Fonteles, Shruti Jain, Xinying Luo, Divya Mereddy, Naveeduddin Mohammed, and Gautam Biswas. 2025. [A theory of adaptive scaffolding for llm-based pedagogical agents](#). *arXiv preprint arXiv:2508.01503*.

Wei Dai, Yi-Shan Tsai, Jionghao Lin, Ahmad Aldino, Hua Jin, Tongguang Li, Dragan Gašević, and Guanliang Chen. 2024. [Assessing the proficiency of large language models in automatic feedback generation: An evaluation study](#). *Computers and Education: Artificial Intelligence*, 7:100299.

David Dinucu-Jianu, Jakub Macina, Nico Daheim, Ido Hakimi, Iryna Gurevych, and Mrinmaya Sachan. 2025. [From problem-solving to teaching problem-solving: Aligning LLMs with pedagogy using reinforcement learning](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 272–292, Suzhou, China. Association for Computational Linguistics.

Juan Escalante, Austin Pack, and Alex Barrett. 2023. [Ai-generated feedback on writing: insights into efficacy and enl student preference](#). *International Journal of Educational Technology in Higher Education*, 20(1):57.

Thomas Palmeira Ferraz, Kartik Mehta, Yu-Hsiang Lin, Haw-Shiuan Chang, Shereen Oraby, Sijia Liu,

974	Luke Mandouit and John Hattie. 2023. Revisiting “the power of feedback” from the perspective of the learner . <i>Learning and Instruction</i> , 84:101718.	1031
975		1032
976		1033
977	Hunter McNichols, Wanyong Feng, Jaewook Lee, Alexander Scarlatos, Digory Smith, Simon Woodhead, and Andrew Lan. 2024a. Automated distractor and feedback generation for math multiple-choice questions via in-context learning . <i>Preprint</i> , arXiv:2308.03234. ArXiv:2308.03234.	1034
978		1035
979		1036
980		1037
981		1038
982		1039
983	Hunter McNichols, Jaewook Lee, Stephen Fancsali, Steve Ritter, and Andrew Lan. 2024b. Can large language models replicate its feedback on open-ended math questions?	1040
984		1041
985		1042
986		1043
987	Jisoo Mok, Ik-hwan Kim, Sangkwon Park, and Sungroh Yoon. 2025. Exploring the potential of LLMs as personalized assistants: Dataset, evaluation, and analysis . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10212–10239, Vienna, Austria. Association for Computational Linguistics.	1044
988		1045
989		1046
990		1047
991		1048
992		1049
993		1050
994	Atharva Naik, Jessica Ruhan Yin, Anusha Kamath, Qianou Ma, Sherry Tongshuang Wu, Charles Murray, Christopher Bogart, Majd Sakr, and Carolyn P Rose. 2024. Generating situated reflection triggers about alternative solution paths: A case study of generative ai for computer-supported collaborative learning . In <i>International Conference on Artificial Intelligence in Education</i> , pages 46–59. Springer.	1051
995		1052
996		1053
997		1054
998		1055
999		1056
1000		1057
1001		1058
1002	Inderjeet Jayakumar Nair, Jiaye Tan, Xiaotian Su, Anne Gere, Xu Wang, and Lu Wang. 2024. Closing the loop: Learning to generate writing feedback via language model simulated student revisions . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 16636–16657.	1059
1003		1060
1004		1061
1005		1062
1006		1063
1007		1064
1008	Tanya Nazaretsky, Hagit Gabbay, and Tanja Käser. 2025. Can students judge like experts? a large-scale study on the pedagogical quality of ai and human personalized formative feedback . <i>Computers and Education: Artificial Intelligence</i> , page 100533.	1065
1009		1066
1010		1067
1011		1068
1012		1069
1013	Sankalan Pal Chowdhury, Vilém Zouhar, and Mrinmaya Sachan. 2024. Autotutor meets large language models: A language model tutor with rich pedagogy and guardrails . In <i>Proceedings of the Eleventh ACM Conference on Learning@ Scale</i> , pages 5–15.	1070
1014		1071
1015		1072
1016		1073
1017		1074
1018	Romain Puech, Jakub Macina, Julia Chatain, Mrinmaya Sachan, and Manu Kapur. 2025. Towards the pedagogical steering of large language models for tutoring: A case study with modeling productive failure . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 26291–26311.	1075
1019		1076
1020		1077
1021		1078
1022		1079
1023		1080
1024	Keyang Qian, Yixin Cheng, Rui Guan, Wei Dai, Flora Jin, Kaixun Yang, Sadia Nawaz, Zachari Swiecki, Guanliang Chen, Lixiang Yan, and 1 others. 2025a. Dean of llm tutors: exploring comprehensive and automated evaluation of llm-generated educational feedback via llm feedback evaluators . <i>arXiv preprint arXiv:2508.05952</i> .	1081
1025		1082
1026		1083
1027		1084
1028		1085
1029		1086
1030		1087
	Keyang Qian, Kaixun Yang, Wei Dai, Flora Jin, Yixin Cheng, Rui Guan, Sadia Nawaz, Zachari Swiecki, Guanliang Chen, Lixiang Yan, and 1 others. 2025b. Scalefeedback: a large-scale dataset of synthetic computer science assignments for llm-generated educational feedback research . <i>arXiv preprint arXiv:2508.05953</i> .	1088
		1089
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300

pages 10364–10379, Albuquerque, New Mexico. Association for Computational Linguistics.

A Surveyed papers by domain

Table 1 presents an overview of all the surveyed papers in this work by their subject domain.

B Mechanisms Supporting Components

B.1 Student state estimation and simulation

Feedback can be guided by an estimate of what the learner *already knows*, leading to decision making over learner state at inference time, or what the learner is *likely to do next*, for feedback policies optimized against predicted responses to feedback.

The student state could be represented by the mastered (or lacking) knowledge bits, like in (Wang et al., 2025b) which uses knowledge tracing for state estimation from dialogue to guide tutor behaviour turn-by-turn. In math tutoring, predicted subsequent success is used as signals to improve tutor actions under sequential constraints (Scarlatos et al., 2025b). In writing, (Nair et al., 2024) uses a student simulator to produce revisions conditioned on candidate feedback and uses revision quality to induce preferences.

Beyond that, student simulators could also be used to facilitate tutoring data creation (Liu et al., 2024a) and system evaluation (Pal Chowdhury et al., 2024). On the other hand, classroom-scale simulation (SimClass) aims to capture multi-party instructional dynamics (Zhang et al., 2025b).

B.2 Scalable feedback aggregation

In classroom settings, qualified experts could be empowered to validate, edit, and reuse aggregated feedback to improve efficiency. The process of *triage, verify, and reuse* decides the common cases to address, the evidence to inspect, and the propagation strategy (Tang et al., 2025b,a), e.g., for large programming course cohorts (Ghoochani et al., 2025), and for live instruction (Wang et al., 2024b).

B.3 Self critique and refinement loops

Refinement loops are general purpose primitives that can be used with feedback systems to improve reliability at inference time under constraints.

Interactive critiquing uses tools, as in (Gou et al., 2024), to ground self-correction in external verification evidence. Agentic educational systems also leverage agents that critique and revise in stages

as a plug-in refinement component within larger workflows (Zhang et al., 2025a).

On the other hand, multi-constraint requirements (e.g. pedagogical principles) can be decomposed by the model’s inherent reasoning/in-context learning into explicit constraints to self-critique violations, and iteratively refine only the parts that fail (Ferraz et al., 2024; Yuan et al., 2024) in single- and multi-agent settings (Jiang et al., 2024; Guo et al., 2024). Iterative revise-and-review loops are also used to improve consistency and constructiveness in reflective writing feedback (Bhojan and Xin, 2025).

C Datasets and Evaluation

Evaluating feedback generation is more challenging than generic text because goodness depends on both **what intervention the system chooses** and **how that intervention is realized in language**. As a result, evaluation should distinguish decision action quality (whether the action is appropriate given the grounding signals) from feedback quality and impact (whether the final feedback is pedagogically valid and improves downstream outcomes).

Across the literature, authors operationalize “good feedback” via observable **evaluation signals**—unit test execution, rubric judgments, revision outcomes, learner-study measures, or instructor edits. When direct learning outcomes are expensive or unavailable, proxies are utilised such as predicted next-turn correctness, simulated revisions, or LLM-as-a-judge rubric satisfaction. Because decision making systems are frequently optimized against these signals, evaluation must clearly state: (i) what is optimized, (ii) what signal supports the claim, and (iii) how trustworthy the signal is (Stamper et al., 2024).

C.1 Datasets

Table 2 in the Appendix provides a summary of the datasets we surveyed, including artifact types, supervision signals, scale, and evaluation protocols.

C.2 Evaluation criteria

We observe four recurring objective families with unique evaluation requirements and threats to validity, but are not necessarily mutually exclusive:

(A) **Feedback pedagogical validity**: feedback that is a valid instructional intervention: content-correct, non-revealing, diagnostic, actionable, and appropriately calibrated in tone and confidence. This naturally pairs with rubric-based evaluations (Scarlatos et al., 2024; Sonkar et al., 2024;

Math	(Baral et al., 2024), (Chudziak and Kostka, 2025), (Dinucu-Jianu et al., 2025), (Liu et al., 2024b), (Liu et al., 2024a), (Jiang et al., 2024), (Puech et al., 2025), (Macina et al., 2025), (McNichols et al., 2024a), (McNichols et al., 2024b), (Pal Chowdhury et al., 2024), (Puech et al., 2025), (Reddig et al., 2025), (Scarlatos et al., 2024), (Scarlatos et al., 2025a), (Scarlatos et al., 2025b), (Tonga et al., 2025) (Wang et al., 2024a), (Wang et al., 2024b), (Weitekamp et al., 2025)
Programming	(Chae et al., 2024), (Ghoochani et al., 2025), (Qian et al., 2025b), (Qian et al., 2025a), (Heickal and Lan, 2024), (Lohr et al., 2025), (Ma et al., 2024), (Naik et al., 2024), (Scholz et al., 2025), (Tang et al., 2025b), (Tang et al., 2025a), (Wang et al., 2024a),
Writing	(Benharrak et al., 2024), (Chamoun et al., 2024), (Göldi et al., 2024), (Chen et al., 2025), (Bhojan and Xin, 2025), (Han et al., 2024), (Nair et al., 2024), (Jovic et al., 2025), (Naik et al., 2024), (Rashkin et al., 2025), (Wang et al., 2025c), (Yuan et al., 2024),
Other domains	(Guo et al., 2024), (Knievel et al., 2025), (Liu et al., 2024b), (Shankar et al., 2025), (Sonkar et al., 2024), (Zhang et al., 2025a), (Stamper et al., 2024), (Nazaretsky et al., 2025), (Seo et al., 2025), (Ferraz et al., 2024), (Yang et al., 2025), (Cohn et al., 2025), (Vanzo et al., 2025), (Zhang et al., 2025b)

Table 1: Surveyed papers by domain

Chen et al., 2025; Han et al., 2024; Liu et al., 2024a; Guo et al., 2024).

(B) Task-level and learning outcomes: feedback that leads to success on the task. In programming, this is often unit test execution success (Chae et al., 2024; Wang et al., 2025b), while in writing, revision quality improvement is usually used as a proxy (Nair et al., 2024), or instead through learner questionnaire (Han et al., 2024). In tutoring, it can be proxied by predicted next-step success under an explicit student model (Scarlatos et al., 2024). Some works measures learning outcome more directly using real participants with pre/post evaluations (Ma et al., 2024; Naik et al., 2024), providing stronger evidence of educational impact but are typically more expensive and harder to standardize across domains.

(C) Classroom efficiency: feedback insights that lead to improvements in throughput and reduce workload in classroom workflows (Tang et al., 2025b,a).

(D) Controllability and stability of pedagogical feedback: Whether a system can (i) follow a requested pedagogical policy (feedback type/criterion/scaffolding level) and (ii) do so consistently across runs. This is evaluated via type adherence in programming feedback (Lohr et al., 2025), coherent multi-step feedback ladders (He-

ickal and Lan, 2024), and repeated-run consistency for rubric-aligned reflective writing feedback (Bhojan and Xin, 2025). In writing assessment, multi-dimensional analytic scoring tests reliable dimension expression (Wang et al., 2025c); criterion-conditioned critique makes control explicit (Yuan et al., 2024); and agentic writing-feedback pipelines test grounded, targeted feedback under planning constraints (Chamoun et al., 2024). In tutoring, policy-level evaluation asks whether scaffolding adapts appropriately and preserves productive failure (Cohn et al., 2025; Puech et al., 2025; Dinucu-Jianu et al., 2025).

C.3 Evaluation signals

Decision making systems are constrained by what signals are observable in the setting. We group measurable evaluation signals into six classes.

Tool-grounded signals: Coding often uses unit tests as an outcome-grounded signal (Chae et al., 2024; Wang et al., 2025b). These are attractive because they are verifiable and scalable, but they primarily measure functional correctness, not pedagogical properties. Tool-grounded evaluation also extends beyond unit tests to tutor/testbed settings where correctness is defined over interaction trajectories: platforms report outcome-based measures (e.g., Success@N / task completion) (Chudziak and

1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350

Kostka, 2025), domain tutors check reasoning over structured artifacts (e.g., circuits) (Knievel et al., 2025), and testbeds/benchmarks score standardized tutor actions (e.g., next-step and bottom-out guidance) (Weitekamp et al., 2025; Macina et al., 2025).

Rubric-based signals: Often feedback is evaluated against rubrics that encode pedagogical constraints such as correctness, non-revelation, actionability, diagnostic value, and tone, like in MCQ diagnostic math feedback (Scarlatos et al., 2024) and Pedagogical Alignment (Sonkar et al., 2024). Rubric-based signals increasingly extend beyond “overall quality.” They include feedback-type adherence in programming (Lohr et al., 2025), coherence of multi-step feedback ladders (Heickal and Lan, 2024), structured criteria for distractor/feedback alignment and pedagogical quality in math/hinting (McNichols et al., 2024a,b; Baral et al., 2024; Tonga et al., 2025), and multi-dimensional analytic writing rubrics that score distinct dimensions (Wang et al., 2025c). As LLM judges become common, studies evaluate judge consistency/accuracy (Seo et al., 2025) and develops validated feedback-evaluator suites (Qian et al., 2025a); preference-based optimization further requires transparent judge setup and calibration (Zhang et al., 2025a). Learner ratings may also diverge from expert judgments due to source credibility effects, motivating designs that disentangle perceived credibility from pedagogical quality (Nazaretsky et al., 2025). When rubrics utilise LLM judges, rater reliability becomes a key requirement. This should be calibrated through agreement with human labellers, consistency checks, or sensitivity analyses. Otherwise, scores could reflect judge bias or surface-form sensitivity rather than genuine improvements (Scarlatos et al., 2024; Tang et al., 2025b).

Revision improvement: Writing feedback is commonly evaluated by its downstream effect on revisions. Revision-based signals increasingly emphasize targeted feedback: scientific writing evaluates whether critiques are text-grounded and actionable for the intended revision (Chamoun et al., 2024); reflective writing measures rubric gains and cross-run consistency, treating iteration as part of the protocol (Bhojan and Xin, 2025); and higher-education deployments pair revision evidence with qualitative analysis of learner trust and adoption (Jovic et al., 2025). (Nair et al., 2024) uses a simulated reviser to generate and score revised drafts for preference training. This

achievement-oriented proxy depends on how well the simulation reflects real learning.

Direct human evidence: Some works measure outcome with human studies. (Ma et al., 2024) reports pre/post-feedback outcomes for debugging skill in novices, while collaborative learning in (Naik et al., 2024) is measured by pre/post reports of interaction evidence. Writing-support and persona-conditioned systems often combine behavioural logs (e.g. feature use) with subjective outcomes (perceived usefulness) (Göldi et al., 2024; Benharrak et al., 2024). Other human studies include: large-course comparisons of peer vs. LLM feedback preferences/actionability (Ghoochani et al., 2025); school RCTs measuring learning/engagement gains from interactive homework tutoring (Vanzo et al., 2025); field studies of human–AI tutoring support in real-time instruction (Wang et al., 2024b); speech assessment pipelines combining ASR quality (e.g., WER) with agreement against human scoring (incl. subgroup analyses) (Shankar et al., 2025); and human judgments validating diagnosis/remediation quality in tutoring (Wang et al., 2024a; Reddig et al., 2025). These evaluations provide stronger evidence from the real world, but are typically small scale and harder to scale across settings.

Instructor-in-the-loop signals: Workflow systems must evaluate both the quality of finalized feedback and the validity of intermediate decision modules. e.g. issue detectors. (Tang et al., 2025b) reports component-level quality alongside human-coded feedback quality and workflow outcomes. (Tang et al., 2025a) evaluates instructor’s feedback review efficiency and edit behaviour, in addition to aggregate measures like misconception coverage.

Benchmark and testbed signals: Standardized benchmarks and interactive testbeds evaluate tutor-like planning policies under shared protocols, e.g., MathTutorBench for open-ended pedagogical capabilities (Macina et al., 2025) and TutorGym for measuring action-level tutoring behaviors across domains (Weitekamp et al., 2025). Scalable evaluation can also rely on synthetic data and simulations, but this increases the need to report synthetic–real fidelity/coverage (Qian et al., 2025b) and to consider how simulation realism may affect conclusions (Zhang et al., 2025b).

1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502

C.4 Validity threats and failure modes

Evaluation conclusions depend on the reliability and validity of the underlying signals. We highlight four categories of validity threats.

Does the metric gauge the real educational construct? Learning science emphasizes evaluating instructional and learner outcomes directly when feasible, and otherwise being explicit about which proxy is optimized and why (Stamper et al., 2024), because passing unit tests might not imply conceptual understanding, rubric satisfaction does not guarantee learning outcome, and revision quality improvement may reflect surface-level edits. A recurring example is that liking or trusting feedback does not necessarily imply learning gains. Large-scale course deployments show that students may prefer LLM feedback and struggle to distinguish it from peer feedback (Ghoochani et al., 2025), while school-based evaluations highlight that stronger evidence comes from designs that include learning outcomes and engagement measures (e.g., pre/post tests and classroom RCTs), not only satisfaction (Vanzo et al., 2025). Similarly, human–AI tutoring support systems may change instructor practices and student experiences in ways that require measuring both pedagogical impact and workflow effects (Wang et al., 2024b).

Are judgments stable and calibrated? Rubric-based evaluation, especially with LLM judges, requires evidence of reliability: agreement with humans, consistency checks, and sensitivity analyses, to protect from judge idiosyncrasies (Scarlatos et al., 2024). LLM judging is not automatically trustworthy: studies evaluate judge stability/accuracy against humans (Seo et al., 2025) and develops validated feedback-evaluator suites trained on human-annotated data (Qian et al., 2025a). Preference “win rates” can be judge/prompt-sensitive, so reporting judge setup and sensitivity is essential for reproducibility (Zhang et al., 2025a); repeated sampling or aggregation can reduce variance when stability is part of the target construct (Bhojan and Xin, 2025).

Is the observed gain caused by better decision making? For systems that combine multiple components, ablations are needed to decide if gains come from better intermediate decisions or simply stronger generation. Component-level evaluation (e.g., verifier calibration, target-selection accuracy, propagation safety) helps isolate what got improved (Tang et al., 2025b,a). This is especially

important for agentic pipelines with multiple interacting steps: targeted ablations can test whether planning improves groundedness/actionability beyond strong prompting (Chamoun et al., 2024), and pedagogy-steering/alignment methods need component-level checks to show gains reflect pedagogical behaviour rather than surface fluency (Puech et al., 2025; Dinucu-Jianu et al., 2025).

Does the gain generalize? The utilisation of simulated students and revisers supports scaling but may not match real classroom distributions. Proxies such as predicted next-turn correctness require validation and should not be interpreted as direct evidence of learning gains (Scarlatos et al., 2025b; Nair et al., 2024; Wang et al., 2025b). Generalization risk is amplified when evidence comes from synthetic or simulated settings: synthetic datasets require reporting fidelity/coverage against real student work (Qian et al., 2025b), and classroom simulations/testbeds depend on the realism of simulated students and tasks (Zhang et al., 2025b; Weitekamp et al., 2025; Macina et al., 2025). Domain-specific tutors add another generalization risk: performance can be brittle if evaluation tasks do not cover the range of authentic student misconceptions (Knievel et al., 2025).

C.5 What to report for comparability

Given diverse evaluation practices, we recommend that feedback generation systems report:

1. *Primary objective/constraints*: whether the goal is pedagogical validity, learning gains, workflow efficiency, or a combination thereof.
2. *Evaluation signals*: whether evaluation uses unit tests, rubrics/judges, revision outcomes, learner-study outcomes, or instructor edits.
3. *Reliability evidence*: report inter-rater agreement for human labellers, and for LLM judge report calibration against humans and stability checks (e.g., repeated-run agreement). Where possible, also report cross-model agreement to diagnose judge sensitivity
4. *Component-level validity*: For verifiers or propagation modules, report calibration/accuracy and safety checks.
5. *Intermediate decision evaluation*: When the system makes explicit intermediate decisions, report accuracy or calibration against reference.

1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550

- 1551 6. *Decision-making design*: specify the intermedi- 1600
 1552 ate decisions the system makes (e.g., feedback 1601
 1553 type, scaffold level, criterion/target selection), 1602
 1554 how these decisions are represented (labels, tem- 1603
 1555 plates, plans), and how they are enforced or 1604
 1556 checked (constraints, verifiers, critique loops). 1605
 1557 If controllability is a goal, report adherence met- 1606
 1558 rics (e.g., type-following rate, ladder coherence, 1607
 1559 criterion coverage) alongside general quality 1608
 1560 scores. 1609
- 1561 7. *Downstream outcome*: e.g. unit-test success, re- 1610
 1562 vision improvement in writing and throughput 1611
 1563 gains for workflows. When feasible, pair sub- 1612
 1564 jective ratings (e.g., perceived usefulness) with 1613
 1565 learning/engagement outcomes using stronger 1614
 1566 study designs (e.g., pre/post tests, randomized 1615
 1567 trials, and classroom field evaluations)
- 1568 8. *Ablations*: to separate the generation effects
 1569 from the contributions of the decision making
 1570 mechanism.
- 1571 9. *Limitations*: include validation evidence and
 1572 scope claims when using simulated stu-
 1573 dents/revisers or proxy models. If train-
 1574 ing/evaluation uses synthetic or simulated data,
 1575 report fidelity/coverage analyses and discuss
 1576 how conclusions might shift under real student
 1577 distributions.

1578 **C.6 Discussion: Evaluation - separating**
 1579 **action quality, feedback quality, and**
 1580 **impact**

1581 A recurring lesson from the surveyed work is
 1582 that feedback evaluation must distinguish (i) *de-*
 1583 *cision/action quality* from (ii) *linguistic quality* and
 1584 (iii) *downstream impact* (§C). Tool-grounded met-
 1585 rics (unit tests) can scale and are verifiable, but do
 1586 not necessarily measure learning or pedagogical
 1587 validity (Chae et al., 2024). Rubric-based judg-
 1588 ments better capture pedagogical constraints (Scar-
 1589 latos et al., 2024; Sonkar et al., 2024), but intro-
 1590 duce rater variance and judge sensitivity, espe-
 1591 cially when LLMs are used as evaluators (Seo
 1592 et al., 2025). Revision improvement provides
 1593 an attractive proxy for writing, yet may conflate
 1594 surface-level edits with genuine learning-relevant
 1595 progress (Nair et al., 2024; Jovic et al., 2025). The
 1596 strongest evidence still comes from human studies
 1597 and deployments that measure learning and engage-
 1598 ment outcomes (Ma et al., 2024; Vanzo et al., 2025),
 1599 but these are expensive and difficult to standardize.

One implication is that evaluation should increas-
 ingly be reported as a *bundle* rather than a single
 number: a primary objective, its signal(s), reliabil-
 ity evidence for those signals, and ablations that iso-
 late whether improvements are attributable to better
 intermediate decisions or simply stronger genera-
 tion. Where “impact” is claimed, stronger designs
 (e.g., pre/post measures or randomized compar-
 isons) are needed to separate learning gains from
 mere revision compliance (Ma et al., 2024; Vanzo
 et al., 2025). Finally, the rise of evaluator-centric
 pipelines (e.g., gatekeeping “Dean”-style evalua-
 tors) suggests that evaluator auditing and calibra-
 tion should be treated as a first-class requirement,
 not an optional appendix (Qian et al., 2025a; Seo
 et al., 2025).

Paper	Task / Domain	Input	Output	Training signal	Size
(Scarlatos et al., 2025b)	Dialogue tutoring (math)	Tutoring dialogue context + student turn (MathDial)	Tutor response	DPO-style preferences favoring tutor turns that elicit correct next student responses; plus GPT-4o distillation to strengthen the base tutor	MathDial: 1,852 (28,125 tutor turns) Test: 361 (5,347 tutor turns) Distill: 6,250; DPO: 11,250 pairs
(Chae et al., 2024)	Code feedback	Problem description + wrong code (with tests in the environment)	NL feedback on how to fix code	Pairwise correct/incorrect feedback labels + unit-test-driven reward (COFFEEVAL) used as supervision for improving feedback	COFFEE: 44,782; 742 problem sets COFFEE-TEST: 180 Avg 35.5 tests/problem
(Wang et al., 2025b)	Coding tutor + verification	Coding task context + dialogue history/turn	Tutor turn + turn-level verifier judgments	Supervised turn-by-turn verifiers used to assess tutoring turns during interactive coding tutoring	EvoCodeBench: 100 tasks
(Nair et al., 2024)	Writing feedback	Student draft (and prompt/context)	Feedback comments	“Closed-loop” learning signal based on LM-simulated revisions conditioned on feedback	363 datapoints (291/72) 873 feedback items
(Sonkar et al., 2024)	Pedagogically scaffolded tutoring	Conversation context (biology tutoring dialogues)	Structured tutor response (incl. pedagogical fields)	Preference learning (DPO/IPO/KTO) on synthetic preference pairs derived from structured tutoring signals (CLASS-style fields)	1,738 conversations total 600 SFT (4,942 QA) 600 LHP (4,921 QA) 450 test (3,701 QA)
(Rashkin et al., 2025)	Story writing feedback (eval dataset)	Story draft (often intentionally corrupted)	Feedback to improve story writing	Human ratings + automatic metrics (dataset is primarily for evaluation/analysis of feedback quality)	Test: 1,300 stories StoryFeedback: ~84k pairs + human-rated subset
(Liu et al., 2024a)	Socratic math tutoring	Math problem + multi-round tutoring dialogue (simulated students)	Socratic tutor turns (questions/hints)	SFT on SocraTeach generated via “Dean-Teacher-Student” pipeline; plus augmentation for specific teaching abilities	SocraTeach: 35k multi-round (=208k single-round equiv.) + 22k extra single-round

Table 2: Summary of publicly available datasets