



Learning Intelligence Engine

Integrate a hybrid student-model combining *knowledge tracing*, *item calibration*, and *reinforcement learning* techniques. Use **Bayesian Knowledge Tracing (BKT)** to model mastery of each topic (skill) over time ¹. BKT is simple and interpretable (AUC $\approx 0.74\text{--}0.80$ on next-question prediction ¹) and requires only moderate data. Complement BKT with **Item Response Theory (IRT)** (2PL/3PL) to calibrate each question's difficulty and discrimination ² ³. For example, calibrate IRT parameters (a , b , c) from past exam data, then estimate a student's ability (θ) from their answer pattern ⁴ ⁵. This allows mapping our score to the NBME scale (e.g. predicting an equivalent Step score ⁶) and selecting items that match student ability (adaptive testing).

Also include an **ELO-style rating** as a lightweight alternative to IRT. Empirical work shows that an ELO model (student & item ratings updated by each question) performs comparably to IRT/BKT while being easier to update online ⁷. We can even cross-check ELO difficulties against IRT from NBME data (Pelánek 2016 found >0.9 correlation of item ordering ⁸). Maintain ELO-based student and question ratings as a quick dynamic measure of ability and difficulty.

For prediction tasks, use **logistic regression and gradient-boosted trees (GBM)** on rich features (student history, timing, question tags, etc.) ⁹ ¹⁰. Logistic regression is a natural baseline (e.g. Performance Factor Analysis uses past successes/failures as features ⁹) and often achieves AUC $\sim 0.75\text{--}0.85$ ¹¹. GBMs (e.g. XGBoost/LightGBM) can capture complex interactions (e.g. timing * past attempts) and have proven high accuracy on benchmark data ¹⁰. Kaggle Riiid contestants found tuned XGBoost reached $\sim 0.78\text{--}0.79$ AUC (compared to ~ 0.82 for deep models) [25†]. Use PyTorch/TensorFlow if exploring deep models, but ensure simpler models as fallback. We should log all student-event data to Snowflake/Redis for model training; train models offline (nightly or batch) and serve predictions via FastAPI/Go services.

In addition, implement **spaced-repetition/forgetting models** for long-term retention. Model each student's *forgetting curve* (e.g. using half-life regression or Anki's SM-2 algorithm) to schedule reviews ¹² ¹³. Duolingo's half-life model increased user retention $\sim 9\text{--}12\%$ ¹³, and Anki users score +11.5% higher on exams ¹⁴. We can fit a per-user HLR (half-life regression) or train a small neural network that predicts recall probability given time since last seen, successes/failures, etc. Integrate this with BKT: use BKT for short-term skill mastery and an HLR model to decide *when* to review older material ¹⁴ ¹⁵. For example, after each answer, update both a BKT state and the memory model; use the memory model's predicted recall to surface questions when their retention drops below a threshold. Also allow the user to rate difficulty/confidence (like Anki's "Again/Good/Easy" feedback ¹⁶) to refine the ease factor. Avoid "overload": detect when too many items are due and throttle the schedule (unlike vanilla Anki) ¹⁷.

To further enhance sequencing, introduce **Contextual Multi-Armed Bandits (MAB)** on top of the learner model. Bandits will choose which topic or question type to present next, balancing review and new material ¹⁸. For instance, treat each "arm" as a content option (next topic Q, review question, flashcard, etc.), with reward = learning gain (e.g. BKT mastery increase or correctness) ¹⁹ ²⁰. A bandit (e.g. Thompson Sampling or UCB) will quickly favor options that improve outcomes. Prior work (e.g. DABSEC) shows bandit-based sequencing can outperform fixed schedules ²¹. Use contextual bandits by feeding current student

state (ability estimate, recent performance) as context features ²⁰. We could even implement an online learning layer (e.g. train a contextual logistic model on the fly for the bandit) ²². Evaluate bandit policies by simulation on historical logs (e.g. replay ASSISTments or Edu records) before live A/B testing ²¹ ²³. For example, simulate sequences on a past cohort log and compare time-to-mastery between different strategies.

Synergies & Implementation: These models should run as a **microservice ensemble**. For each student action, update BKT/IRT/ELO states in a Python/Go service and log outcomes. The bandit service queries the current state and returns the next best action. Batch jobs (e.g. nightly) will train logistic/GBM on Snowflake data and export parameters (or m2cgen code) for serving. Use Redis to cache hot features (recent BKT scores) for low-latency. Tie all learner-model services via events (OpenTelemetry traces) to monitor performance. Evaluate models via offline AUC accuracy on held-out data and live A/B tests measuring actual learning (e.g. improvement on practice exams). Use the Riiid/EdNet dataset (public Kaggle data) or ASSISTments for initial validation ²⁴. Architecturally, each algorithm can be its own containerized service, allowing horizontal scaling (e.g. large GPU nodes for heavy retraining of PyTorch models).

CMS and Content Ingestion

Adopt a **component content management system (CCMS)** with structured authoring and rich metadata. Each question/item, its explanations, and related content (images, references) should be stored as modular components (e.g. Markdown or XML fragments) following a medical taxonomy. Use industry standards (e.g. IMS QTI for assessment items) so content is portable ²⁵. Enforce a **schema/taxonomy** for questions: tag each item with concepts (mapped to our Neo4j concept graph nodes), organ system, difficulty level, Bloom's taxonomy level, etc. This advanced tagging enables powerful search and analytics (linking Qs to prerequisite concepts, auditing difficulty).

Build an **editorial workflow** with roles and version control. For example, integrate content with Git or a headless CMS (like Strapi/Sanity) that supports branching and pull-requests. Implement review queues: when a question is authored or edited, require approval by a subject expert. Use automated checks (e.g. grammar/spell-checkers for medical terms) in the pipeline. Track full **version history and audit logs** for every change ²⁶ ²⁵. (Moodle's question bank highlights this: it provides versioning for Qs, showing change history and usage ²⁶.) Allow editors to comment on or mark questions as "draft" or "approved" ²⁶.

Integrate a **review and feedback loop** into the CMS. Each time students answer a question, collect metadata: empirical difficulty (proportion correct), discrimination index, and student feedback (e.g. flag "ambiguous" or rate clarity). Store these in the CMS database linked to the question. Like TAO's authoring platform, capture item statistics and quality hints ²⁵. If a question's performance falls (e.g. everyone gets it wrong), flag it for editorial review. Use BI dashboards to surface questions with low discrimination or high dropout.

For content creation, provide authors with rich tooling: a React-based editor (with LaTeX/math support, image upload, and links to references). Support **component reuse**: for example, store common definitions or diagrams as separate modules that can be embedded. This follows structured authoring best practices ²⁷. Templates and style guides (for headings, citations, formulas) ensure consistency and compliance with medical standards ²⁷.

Overall, the CMS should feel “medical-grade”: content is peer-reviewed, quality-assured, and traceable. We may borrow ideas from documentation CCMS: use taxonomies (SNOMED/MeSH) for tagging and enforce an approval workflow [27](#) [28](#). Enable export/import (e.g. to Anki decks or other systems), and ensure a full audit trail (who changed what and when) [27](#) [26](#).

Analytics & Prediction

Implement a suite of dashboards for learners, instructors, and admins. Use modern JS libraries (D3.js is in our stack, but also consider simpler charting like Chart.js or Plotly) to build interactive charts. Follow best practices: highlight key metrics (progress bars, masteries achieved) with tooltips/drill-down [29](#) [30](#), and use meaningful color-coding (e.g. green/yellow/red difficulty) as seen in Amboss [31](#).



Figure: Example analytics dashboard (students vs. topics).

For cohort analytics, allow comparing a student to peers or historical groups. For example, show percentile rank, heatmaps of cohort performance by subject, or “learning trajectory” charts. Use small multiples (one chart per subject) or boxplots to show distribution of scores across peers. Leverage libraries like Vega-Lite or Recharts for clear dashboards.

As an improvement, build a **rank simulation** feature: given current performance and model-estimated growth, simulate a student’s likely final exam rank or score. This could use our IRT ability θ and NBME-equated models. For instance, if the student’s current θ predicts a Step1 score of 240 (based on NBME mapping [6](#)), show a projected percentile. Running “what-if” scenarios (if effort increases by 20%, estimate rank change) can motivate students.

Ensure all visualizations minimize cognitive load [29](#) [30](#): keep dashboards uncluttered, use consistent layouts, and present data progressively (summary first, details on click). For example, the main dashboard

might show “Your Overall Mastery: 75%” with a breakdown bar-chart by subject, then clicking a subject reveals topic-level charts.

Search & Exploration

Enhance search by combining Elasticsearch keyword queries with semantic vector search. Index all questions and content in Elasticsearch; use full-text fields, metadata filters (subject, difficulty), and tagging. For semantic discovery, compute embeddings of question text or concepts (using an OpenAI/HuggingFace model or the built-in ES `semantic_text` field ³²). Then perform a k-NN vector search to find semantically similar questions or relevant resources. Elastic’s `dense_vector` search retrieves by meaning (not just exact terms) ³² ³³, ideal for “questions on kidney stones” even if phrased differently.

Also leverage the **Neo4j concept graph** for exploration. When a student searches or reviews a topic, use Cypher queries on Neo4j to suggest prerequisite or related concepts. For example, if a user looks up “atrial fibrillation,” traverse graph edges to show underlying topics (“electrophysiology basics → atrial flutter → atrial fibrillation”). This graph-aware path planning helps students see a logical learning sequence ³⁴. In question discovery, query Neo4j to find all questions linked to concepts along a path. For advanced use, combine vector + graph: embed user query to get candidates via ES, then rank/filter by concept graph distance. This **GraphRAG** approach (retrieval-augmented generation with a knowledge graph) yields semantically and topically relevant results ³⁵.

Offer an interactive concept map UI: when studying a topic, display its neighbors (prerequisites and outcomes). Neo4j can serve this data; use D3 or a library like Cytoscape.js to render the graph. This helps students plan (e.g. “I need to review ‘renal physiology’ before tackling ‘acid-base disorders’”).

Security & Integrity

Protect exam content and user data with best practices. Use **end-to-end encryption** (TLS everywhere). Sensitive payloads (e.g. question data, answer submissions) should be encrypted-in-transit and ideally encrypted-at-rest. For high-stakes tests, consider encrypting exam bundles so they cannot be read without a secure key when the exam starts.

Implement **device fingerprinting** to deter fraud. Use a service like FingerprintJS: it collects browser/device signals and returns a device ID along with risk signals. Enable their *sealed client results* mode, which encrypts the fingerprint payload so it can’t be tampered with client-side ³⁶. On each login, record the fingerprint and detect anomalies (new device, location change). Flag suspicious patterns (VPN use, rapid switching).

Enforce **rate limiting and session controls**: throttle login attempts, answer submissions, and API calls to prevent brute force or scraping. Use CAPTCHA on suspicious login flows. If hosting real-time tests, implement lockdown measures (block printing, clipboard).

Architect a **Zero Trust** environment: do not assume internal networks are safe. Require authentication/authorization on every microservice call (e.g. validate JWT tokens or use mutual-TLS between services). Store secrets in a vault (AWS Secrets Manager/Azure Key Vault). Monitor and log all access events.

Lastly, pursue security compliance (e.g. ISO 27001, SOC-2) as needed for a high-stakes platform. Regularly perform penetration testing and vulnerability scans.

UX/UI and Visualization

Design interfaces that **minimize cognitive load** ²⁹ ³⁰. For test-taking mode, present one question at a time on a clean layout: large fonts, high-contrast text, and a visible countdown timer. Use color-coding sparingly (e.g. green for correct, red for flagged weakness). As Amboss does, include visual cues like a “preparedness meter” or color highlights on weak topics ³¹. Provide a “tutor mode” where hints or partial answers are revealed progressively, and an “exam mode” where no hints are shown (like UWorld’s Attending Mode ³¹).

For dashboards and review, use modern JS frameworks (React/Next) with components (Tailwind CSS for consistency). Use animation libraries judiciously: GSAP (GreenSock) can add smooth transitions (e.g. progress bar fills, card flips for flashcards) but ensure animations do not distract or add load ²⁹. Use D3.js for custom visualizations (e.g. a skill tree graph or knowledge map). For simpler charts (bar, line, pie), Chart.js or Nivo can speed development.

Implement interactive visual aids: for example, an “exam map” that shows which topics remain unmastered, or a calendar heatmap of study streaks. Follow UX patterns like “depth-on-demand” (collapsible explanations) ³⁷ and immediate feedback on actions. Ensure responsiveness (support tablets/phones) and accessibility (ARIA labels, keyboard navigation). Collect UX analytics (via Hotjar or similar) to find confusing elements and iterate.

Mobile Experience

Build the React Native app **offline-first**: bundle a subset of questions and content in a local database (e.g. SQLite or WatermelonDB). Follow the “local-first” principle ³⁸: use the local DB as the source of truth, present cached data immediately, and sync changes when the network is available ³⁸. For example, if a student starts a practice session offline, log their answers locally and upload when back online. Use service workers or background sync (in a PWA context) to pre-fetch next sessions.

Optimize the UI for *time-pressure exams*: large tappable areas, one-question-per-screen, and minimal navigation. Lock orientation if needed. Show a persistent timer and question count (e.g. “Q4 of 40”). Allow swiping to skip or flag a question for later review. Provide offline alerts (e.g. “Network lost; your progress will sync when restored”). Offer lightweight data visualization for mobile (small sparklines or progress bars).

For mobile charts, consider libraries like Victory Native or React Native Charts, but keep visuals simple to avoid slowing the app. Ensure the app degrades gracefully on low-end devices (avoid heavy vector graphics). Provide an optional “low-data mode” (e.g. skip loading explanation images).

Finally, integrate device features: push notifications for spaced-repetition reminders, and a quick “resume where you left off” link. Maintain user session even if the app is backgrounded, to prevent accidental timeouts.

Testing & Observability

Set up comprehensive testing pipelines. Use **property-based testing** (Hypothesis in Python, QuickCheck) to validate invariants (e.g. shuffling answers doesn't change correctness, no duplicate IDs in Qs). For API endpoints, use contract tests (Swagger-generated tests). Implement **load testing** (Locust, k6 or JMeter) to simulate thousands of concurrent students. For example, simulate a mock cohort taking an NBME-style test to measure response times and resource usage. Perform regular chaos experiments in staging: e.g. kill random pods or inject latency (using Chaos Mesh or Gremlin) while monitoring the system.

Instrument everything with OpenTelemetry and Prometheus. Tag and trace all user actions end-to-end. Use Grafana to build dashboards of key metrics (CPU/memory, request latencies, error rates, queue depths). As a best practice, set up alerts (e.g. high error rate triggers PagerDuty). During chaos tests or spikes, observe traces to pinpoint bottlenecks. For example, using an APM like SkyWalking or Jaeger will let us see how injecting a node failure cascades through service calls ³⁹. Continuously analyze logs (ELK or similar) to catch new errors.

In addition, include analytics-specific observability: monitor model accuracy drift (if AUC drops for live data), and log data quality issues (e.g. missing tags on content, NaNs in response logs). Treat ML model inference as part of observability: log prediction inputs and outputs (anonymized) to spot skew.

Parallelization & DevOps

Exploit parallelism wherever possible. All microservices (user service, question service, analytics) are decoupled and can scale independently. Model training tasks can run in parallel: e.g. train separate GBMs for each subject or run hyperparameter searches concurrently. Use Kubernetes to run GPU/PyTorch jobs parallelized (Distributed Data Parallel or Ray for training). Batch ETL (extract logs from Redis to Snowflake, train models) should use Spark or Dask to parallelize across clusters.

For CI/CD, containerize all components and use GitHub Actions or Azure Pipelines to build on each commit. Use **blue-green deployments**: keep the old version running while testing the new one, then switch traffic once stable ⁴⁰. (Azure docs confirm using blue/green or canary for non-breaking changes ⁴¹.) Maintain separate environments (dev/test/prod) with IaC (Terraform) for consistent configs. Automate database migrations (e.g. Postgres schema via Flyway). Implement canary releases for major features (e.g. rollout new bandit algorithm to 10% of students first).

Leverage Kubernetes features: use readiness/liveness probes, pod anti-affinity for high availability, and resource quotas to prevent noisy neighbors. Use autoscaling (HPA/VPA) for stateless services. Keep staging as a mirror of prod (ideally with anonymized prod data) to test load/release candidates.

Roadmap Completion Plan

Develop iteratively across teams with clear milestones:

1. **Core engine foundation:** Finalize user and content schemas (Postgres/Redis/Neo4j). Build basic services: user auth, question CRUD, tagging. Deploy CI/CD pipeline and Kubernetes cluster.
2. **Content ingestion & CMS:** Launch the CMS/editor interface. Migrate initial question bank. Implement versioning, tagging, and review workflow. Integrate with Git or headless CMS.
3. **Initial intelligence models:** Implement BKT and IRT modules with open datasets. Deploy as background service; validate with historic logs. Train initial GBM on legacy data. Measure baseline AUC.
4. **Feedback and retention:** Add forgetting model (e.g. HLR) and flashcard mode (Anki-like). Hook into CMS tagging for spaced review schedules. A/B test retention vs static review.
5. **Adaptive sequencing:** Deploy bandit engine. Run simulations on logs for tuning. Launch live pilot where some students get bandit-driven sequencing vs control. Monitor learning gains.
6. **Analytics UI:** Build initial dashboards (user progress, cohort stats). Populate with demo data. Refine with user feedback. Optimize charts (use [79]example as prototype).
7. **Search & graph features:** Populate Elasticsearch index and Neo4j graph. Implement semantic question search and prerequisite path finder. Test concept recommendations in UX.
8. **Security hardening:** Integrate fingerprinting JS, implement “sealed” payload flow ³⁶. Add rate-limit middleware. Conduct pen test and fix issues.
9. **UX polish:** Iterate on UI per cognitive load heuristics ²⁹. Add animations (GSAP) for smoothness, ensure mobile responsiveness. Conduct usability tests.
10. **Load/chaos testing:** Perform final scalability tests. Tune K8s scaling rules. Finalize SLOs/alerts in Prometheus+Grafana.
11. **Launch V2:** Do a blue-green release. Migrate content and user base. Announce features (e.g. predicted scores, flashcards). Provide training docs.

At each step, include unit/regression tests and performance benchmarks. Avoid shortcuts: e.g. do *not* skip rigorous tagging or model validation, as that can undermine adaptivity. For example, launching without version control for content could later cause errors; or skipping offline mobile support will frustrate users. Track metrics continuously: do not assume “if it builds, it works” – measure learning outcomes and iterate.

This blueprint, backed by academic insights and proven practices ¹ ¹³ ³⁶, will guide the team to a robust, scalable V2 platform.

Sources: Published adaptive learning research, industry best practices, and the provided MBBS engine reference ¹ ¹⁸ ² ⁷ ³² ²⁹ ¹⁴ ³⁶. Tables and images illustrate architecture and data.

¹ ² ³ ⁴ ⁵ ⁶ ⁷ ⁸ ⁹ ¹⁰ ¹¹ ¹² ¹³ ¹⁴ ¹⁵ ¹⁶ ¹⁷ ¹⁸ ¹⁹ ²⁰ ²¹ ²² ²³ ²⁴ ³¹ ³⁷ Personalized

Learning Engine for MBBS MCQ Exams.pdf

file:///file_000000063807206ac31025c2319df93

²⁵ Test Creation Software | Online Exam Builder| TAO Testing

<https://www.taotesting.com/authoring/>

26 Question bank - MoodleDocs

https://docs.moodle.org/4x/sv/Question_bank

27 28 Structured content for medical writing - RWS

<https://www.rws.com/content-management/tridion/medical-devices/medical-writing/>

29 30 User interface design in mobile learning applications: Developing and evaluating a questionnaire for measuring learners' extraneous cognitive load - PMC

<https://pmc.ncbi.nlm.nih.gov/articles/PMC11422584/>

32 33 Vector search in Elasticsearch | Elastic Docs

<https://www.elastic.co/docs/solutions/search/vector>

34 jcst.ict.ac.cn

<https://jcst.ict.ac.cn/fileup/1000-9000/PDF/2021-5-18-0328.pdf>

35 Using a Knowledge Graph to implement a RAG application

<https://neo4j.com/blog/developer/rag-tutorial/>

36 Sealed Client Results

<https://dev.fingerprint.com/docs/sealed-client-results>

38 Build an offline-first app | App architecture | Android Developers

<https://developer.android.com/topic/architecture/data-layer/offline-first>

39 Chaos Mesh + SkyWalking: Better Observability for Chaos Engineering | Chaos Mesh

<https://chaos-mesh.org/blog/better-observability-for-chaos-engineering/>

40 41 CI/CD for microservices - Azure Architecture Center | Microsoft Learn

<https://learn.microsoft.com/en-us/azure/architecture/microservices/ci-cd>