

# FDA AI Credibility Assessment in Drug Development

By Adrien Laurent, CEO at IntuitionLabs • 3/4/2026 • 55 min read

fda ai framework

ai credibility assessment

drug development

ai model validation

regulatory compliance

context of use

machine learning

risk-based validation



# Executive Summary

The U.S. Food and Drug Administration (FDA) has recently introduced a **risk-based credibility assessment framework** for AI models used in drug and biological product development <sup>(1)</sup> [regulations.justia.com](#) <sup>(2)</sup> [emergingaihub.com](#). In January 2025 the FDA released draft guidance (*Considerations for the Use of Artificial Intelligence to Support Regulatory Decision-Making for Drug and Biological Products*) that outlines **seven steps** sponsors should follow to ensure any AI/ML model's output is trustworthy for its intended regulatory context <sup>(1)</sup> [regulations.justia.com](#) <sup>(2)</sup> [emergingaihub.com](#). The cornerstone of the framework is *context of use (COU)* – i.e. the precise question or decision the model supports – and *credibility evidence*, meaning the data, tests, and documentation needed to establish trust in the model's performance for that COU <sup>(1)</sup> [regulations.justia.com](#) <sup>(2)</sup> [emergingaihub.com](#).

A **risk-based approach** pervades the guidance. Models with high influence on decisions or potentially severe consequences (e.g. patient safety or product quality) require **rigorous validation** and extensive documentation, whereas lower-risk uses may need only basic supporting evidence <sup>(3)</sup> [www.jdsupra.com](#) <sup>(4)</sup> [emergingaihub.com](#). High-stakes applications (e.g. patient selection, dose assignment, or batch-release decisions) will generally demand disclosure of the full model architecture, training data, test results, and uncertainty analysis <sup>(5)</sup> [www.jdsupra.com](#) <sup>(6)</sup> [pmc.ncbi.nlm.nih.gov](#). Conversely, purely supportive uses or discovery-stage models (outside FDA scope) may avoid deep scrutiny <sup>(7)</sup> [regulations.justia.com](#) <sup>(8)</sup> [emergingaihub.com](#).

Key components of evaluation include ensuring **high-quality, "fit-for-use" data**, robust **model validation** on independent datasets, and comprehensive performance metrics (accuracy, confidence intervals, bias measures, etc.) <sup>(6)</sup> [pmc.ncbi.nlm.nih.gov](#) <sup>(9)</sup> [emergingaihub.com](#). The guidance emphasizes documenting all activities: defining the regulatory question (Step 1), specifying COU (Step 2), performing a risk analysis (Step 3), planning and executing validation studies (Steps 4–5), and compiling results and deviations (Step 6) <sup>(1)</sup> [regulations.justia.com](#) <sup>(10)</sup> [emergingaihub.com](#). Finally, sponsors must judge whether the model is "adequately credible" for its COU (Step 7), considering residual uncertainty and planning for continuous monitoring if inputs or operations change <sup>(11)</sup> [emergingaihub.com](#) <sup>(12)</sup> [www.sciencedirect.com](#). Throughout, transparency is stressed: FDA expects clear reporting of model development, testing, and changes so reviewers can **independently evaluate** the AI evidence without rerunning analyses <sup>(13)</sup> [emergingaihub.com](#) <sup>(14)</sup> [www.foley.com](#).

This report provides an in-depth treatment of the FDA's AI credibility framework. We first present historical context on AI/ML in drug development and prior regulatory initiatives. We then unpack each of the seven steps with detailed analysis and expert commentary, covering how to define questions, set COUs, assess risk (influence × consequence), and plan validation. Sections on **data quality, model validation, and performance evaluation** delve into the technical evidence (data curation, bias analysis, uncertainty quantification) needed for credibility. We also discuss **documentation practices**, including what sponsors should include in submissions (data provenance, code descriptions, **audit trails, Part 11 compliance**) and how to balance transparency with intellectual property concerns <sup>(15)</sup> [www.foley.com](#) <sup>(16)</sup> [www.foley.com](#).

Multiple real-world examples illustrate these principles. For instance, a case study in bioprocess simulation describes using an AI "hybrid" model to optimize antibody manufacturing parameters; the sponsor tied the simulation model to a specific critical process parameter (CPP) decision and collected historical batch data and quality tests to validate the model <sup>(17)</sup> [www.ark-biotech.com](#) <sup>(18)</sup> [www.sciencedirect.com](#). In another example, a machine-learning model identified which COVID-19 patients would benefit most from the drug anakinra; this predictive model was used under an FDA emergency use authorization and was supported by retrospective clinical data analysis, demonstrating how AI can guide **patient stratification** <sup>(19)</sup> [pmc.ncbi.nlm.nih.gov](#). We also review AI applications such as EMA-qualified digital pathology tools for NASH liver biopsies and AI-driven real-world evidence platforms, showing how regulatory agencies are already encountering AI in submissions <sup>(20)</sup> [pmc.ncbi.nlm.nih.gov](#) <sup>(6)</sup> [pmc.ncbi.nlm.nih.gov](#).

Finally, we explore broader implications and future directions. The FDA's emphasis on transparency and rigor is expected to spur innovations – for example, tools for explainable AI, automated report generation, and real-time model monitoring under GMP conditions (<sup>[21]</sup> [www.foley.com](http://www.foley.com)) (<sup>[18]</sup> [www.sciencedirect.com](http://www.sciencedirect.com)). Sponsors should prepare by engaging FDA early and aligning their validation plans with the model's risk level. We compare FDA's approach to emerging global guidelines (e.g. the EMA's reflection paper on AI in the pharmaceutical lifecycle) and note remaining challenges in harmonization ([www.ema.europa.eu](http://www.ema.europa.eu)) (<sup>[22]</sup> [emergingaihub.com](http://emergingaihub.com)). Ultimately, the FDA's framework lays out a path for building *trustworthy AI* in drug development – requiring clear COUs, strong evidence, and ongoing oversight – so that AI can deliver its promise of innovation while upholding patient safety and product quality (<sup>[1]</sup> [regulations.justia.com](http://regulations.justia.com)) (<sup>[2]</sup> [emergingaihub.com](http://emergingaihub.com)).

► **Table: Seven Steps of FDA's AI Model Credibility Framework**

## Introduction and Background

Artificial intelligence (AI) and machine learning (ML) are **transforming drug development**, spanning tasks from target identification to clinical trial design and manufacturing optimization (<sup>[26]</sup> [intuitionlabs.ai](http://intuitionlabs.ai)). Sponsors increasingly use AI/ML to predict patient outcomes, identify trial cohorts, analyze biomarkers, and optimize bioprocesses. For example, ML models can integrate large real-world datasets to find patient subgroups or simulate “digital twin” trials (<sup>[26]</sup> [intuitionlabs.ai](http://intuitionlabs.ai)) (<sup>[27]</sup> [pmc.ncbi.nlm.nih.gov](http://pmc.ncbi.nlm.nih.gov)). FDA Commissioner Califf has noted AI's “transformative potential to advance clinical research” when proper safeguards are in place (<sup>[28]</sup> [www.fda.gov](http://www.fda.gov)).

However, these advanced models also introduce new challenges. Many modern AI systems (e.g. deep neural networks or ensemble algorithms) are essentially “black boxes” whose internal logic is not easily interpretable (<sup>[29]</sup> [academic.oup.com](http://academic.oup.com)). This opacity raises concerns: biases or errors in training data may be amplified, model predictions may not generalize across patient groups, and complex models may be difficult for regulators to verify (<sup>[29]</sup> [academic.oup.com](http://academic.oup.com)) (<sup>[6]</sup> [pmc.ncbi.nlm.nih.gov](http://pmc.ncbi.nlm.nih.gov)). Given that drug development decisions directly affect patient safety and product quality, regulators must ensure AI tools are *reliable and valid* for their intended use.

The FDA has long encouraged innovation but has also needed to clarify how AI/ML fits into existing regulatory pathways. Over the past decade, the agency has seen a surge in submissions involving AI “components” (<sup>[30]</sup> [www.fda.gov](http://www.fda.gov)) (<sup>[23]</sup> [regulations.justia.com](http://regulations.justia.com)). For instance, since 2016 FDA centers have reviewed on the order of 300–1,000 drug- or biologics-related submissions incorporating AI/ML methods (<sup>[30]</sup> [www.fda.gov](http://www.fda.gov)) (<sup>[31]</sup> [emergingaihub.com](http://emergingaihub.com)). Common applications have included predictive models for clinical trial enrollment, digital biomarkers, and advanced analytics for manufacturing and pharmacovigilance. In tandem, stakeholders (industry, academia, patients) have pushed FDA for guidance on how AI should be evaluated. FDA has organized workshops (e.g. with Duke's Margolis Center in 2022) and published discussion documents in 2023 (soliciting **800+ comments**) on AI in drug development and manufacturing (<sup>[32]</sup> [www.fda.gov](http://www.fda.gov)) (<sup>[23]</sup> [regulations.justia.com](http://regulations.justia.com)).

In January 2025, the FDA issued its first draft guidance focused on *drug development AI* (complementing a 2023 guidance on AI-enabled medical devices) (<sup>[33]</sup> [www.fda.gov](http://www.fda.gov)) (<sup>[34]</sup> [intuitionlabs.ai](http://intuitionlabs.ai)). Titled “*Considerations for the Use of Artificial Intelligence to Support Regulatory Decision-Making for Drug and Biological Products*”, the draft guidance provides recommendations on how sponsors should plan, validate, and document AI models across the drug lifecycle (nonclinical, clinical, manufacturing) when those models produce data or analyses that will inform regulatory decisions (<sup>[35]</sup> [regulations.justia.com](http://regulations.justia.com)) (<sup>[36]</sup> [www.dlapiper.com](http://www.dlapiper.com)). The FDA emphasizes that *context of use* (COU) is key: each model must have a narrowly defined role addressing a specific safety, efficacy, or quality question (<sup>[35]</sup> [regulations.justia.com](http://regulations.justia.com)) (<sup>[2]</sup> [emergingaihub.com](http://emergingaihub.com)). Crucially, the guidance introduces a **seven-step credibility assessment** that sponsors can follow to demonstrate trust in a model's output. This framework upholds FDA's traditional standards for evidence, framing them as *credibility evidence* for AI–ML models (<sup>[35]</sup> [regulations.justia.com](http://regulations.justia.com)) (<sup>[2]</sup> [emergingaihub.com](http://emergingaihub.com)).

Notably, the draft guidance explicitly **excludes** early drug discovery (e.g. target ID, lead optimization) and purely internal uses (workflow tools, administrative tasks) that do not impact patient safety or product quality (<sup>[7]</sup> [regulations.justia.com](http://regulations.justia.com)) (<sup>[8]</sup> [emergingaihub.com](http://emergingaihub.com)). In other words, it applies only when “AI-generated output is intended to support a regulatory

decision” about a drug’s quality, safety, or efficacy (<sup>[35]</sup> regulations.justia.com). For example, an ML model that flags adverse events or predicts trial outcomes falls within scope, whereas an AI tool used only to design molecules (before any submission) does not (<sup>[7]</sup> regulations.justia.com) (<sup>[8]</sup> emergingaihub.com). By focusing on decision-impacting AI, the guidance seeks to ensure that high-stakes uses of AI are subject to rigorous review.

The FDA’s draft emphasizes **transparency and documentation**. It describes “credibility” as “*trust, established through the collection of credibility evidence, in the performance of an AI model for a particular [COU]*” (<sup>[1]</sup> regulations.justia.com). Credibility evidence can be any data or analysis that supports the model’s validity (training data details, validation results, etc.). The agency encourages early communication with FDA staff to agree on appropriate credibility activities based on model risk (<sup>[37]</sup> regulations.justia.com) (<sup>[23]</sup> regulations.justia.com). As the draft states, from COU to documentation, the proposed recommendations “will generally be tailored to the specific COU and will depend on model risk” (<sup>[23]</sup> regulations.justia.com).

This research report draws on the FDA’s text and public statements (<sup>[30]</sup> www.fda.gov) (<sup>[35]</sup> regulations.justia.com), expert commentaries (<sup>[38]</sup> www.foley.com) (<sup>[36]</sup> www.dlapiper.com), academic analyses (<sup>[29]</sup> academic.oup.com) (<sup>[6]</sup> pmc.ncbi.nlm.nih.gov), and case examples to unpack the FDA’s AI guidance. We detail each of the seven credibility steps, discuss data and validation considerations (e.g. metrics, bias, uncertainty), and review how sponsors should document their AI work. We also survey multiple perspectives—from legal experts, industry practitioners, and international regulators—on opportunities and challenges. Real-world examples illustrate key points: for instance, simulation models in drug manufacturing (<sup>[17]</sup> www.ark-biotech.com) (<sup>[18]</sup> www.sciencedirect.com) and ML-driven trial designs in practice (<sup>[19]</sup> pmc.ncbi.nlm.nih.gov) (<sup>[20]</sup> pmc.ncbi.nlm.nih.gov). By the end, readers will have a comprehensive view of how to evaluate an AI model’s credibility in the context of drug development, how to assemble the required evidence, and what this evolving regulatory landscape implies for the innovation of AI in pharmaceuticals.

## FDA’s AI Credibility Framework: Defining Context and Risk

The heart of the FDA’s draft guidance is a **risk-based credibility framework** for AI models that influence drug development decisions (<sup>[35]</sup> regulations.justia.com) (<sup>[36]</sup> www.dlapiper.com). This framework revolves around two foundational concepts:

- **Context of Use (COU).** The COU precisely defines *how* and *where* the AI model will be applied to a specific decision. It includes the model’s scope, inputs and outputs, and its role in the workflow (e.g. a screening tool vs. an automated classifier). For example, a COU might be “using an ML model on EHR data to predict which trial candidates meet inclusion criteria, to inform investigator referrals” (flagging with human review), versus “AI automatically assigns batch-release clearance in manufacturing without human check” (<sup>[10]</sup> emergingaihub.com) (<sup>[3]</sup> www.jdsupra.com). Clarifying COU ensures the model is evaluated *only for its intended task*, not generalized beyond it (<sup>[1]</sup> regulations.justia.com) (<sup>[10]</sup> emergingaihub.com).
- **Credibility Evidence.** For each COU, sponsors must gather evidence that the model can be trusted. The guidance defines credibility as “trust...established through the collection of credibility evidence, in the performance of an AI model for a particular COU” (<sup>[1]</sup> regulations.justia.com). In practice, this means documenting the model design, training, and validation so that reviewers can independently assess its reliability. Credibility evidence can include training data descriptions, performance metrics (accuracy, AUC, calibration), statistical uncertainty (confidence intervals, error rates), validation on hold-out data, results of robustness tests, etc.—any data relevant to making the model’s output believable for the COU (<sup>[6]</sup> pmc.ncbi.nlm.nih.gov) (<sup>[9]</sup> emergingaihub.com).

The FDA’s credibility approach is explicitly **risk-based**. The agency asks sponsors to consider two axes for model risk: **Model Influence** (how much the AI output drives a decision) and **Decision Consequence** (the seriousness of a wrong decision) (<sup>[4]</sup> emergingaihub.com). A fully automatic decision with no human oversight is “high influence”; a slight aid that a clinician can override is “low influence.” A decision about patient dosing or release of unsafe product is “high consequence”; a purely organizational prediction (e.g. internal resource planning) is “low consequence.” The guidance

notes that “models making final determinations without human oversight generally carry higher risk and thus demand stronger evidence” <sup>(39)</sup> [intuitionlabs.ai](#)). FDA illustrates this with a hypothetical: an AI model classifying trial participants as low-risk so they can be sent home versus inpatient monitored. If the model errs and misclassifies a dangerously high-risk patient as low-risk, the consequences (missed adverse events) are severe, yielding a very high risk scenario <sup>(4)</sup> [emergingaihub.com](#)).

In practical terms, most uses within the guidance’s scope will tend toward high-risk. As noted by Foley LLP, AI used for trial management or manufacturing likely *influences* critical decisions, so sponsors should plan to submit extensive details about the AI model <sup>(5)</sup> [www.jdsupra.com](#)). Indeed, the draft guidance emphasizes that high-risk models are expected to include **comprehensive disclosures**: the model’s architecture, full description of training and source data, algorithms, validation methods, and performance statistics <sup>(5)</sup> [www.jdsupra.com](#) <sup>(6)</sup> [pmc.ncbi.nlm.nih.gov](#)). Lower-risk models might allow pared-down submission (e.g. summary performance). But in either case, the principle is the same: the higher the potential impact, the more rigor needed in building and documenting trust.

Below we outline the first key steps in setting up an AI credibility strategy: defining the question and context, and assessing risk. Later sections delve into the validation and documentation needed to address that risk.

## Step 1: Define the Question of Interest

The **first step** is to articulate the precise “question of interest” the AI model is meant to answer <sup>(2)</sup> [emergingaihub.com](#) <sup>(36)</sup> [www.dlapiper.com](#)). This is essentially a problem statement tied to the FDA’s review goals (safety, efficacy, quality). It forces sponsors to ground their AI use-case in a concrete regulatory decision rather than a vague business goal. For instance:

- In a clinical trial setting, the question might be: “Which patients in trial XYZ are eligible for outpatient monitoring after dosing Drug A, without compromising safety?” (This implicates patient safety, since missing an adverse event could be serious.)
- In manufacturing, a question could be: “Can batch Y, produced under process conditions Z, be released as meeting fill-volume specifications?” (This affects product quality.)

Defining the question is critical because it determines both the *inputs/outputs* of the model and the *acceptance criteria*. As Foley LLP observes, examples include using AI for *trial inclusion criteria, patient risk stratification, or endpoint adjudication* <sup>(40)</sup> [intuitionlabs.ai](#) <sup>(36)</sup> [www.dlapiper.com](#)). Importantly, sponsors should specify whether the AI output is **informing** a decision (with subsequent human review) or making the decision **autonomously**. An AI model that “suggests” labels for a human reviewer is lower risk than one that “decides” on its own. The guidance explicitly notes that higher autonomy increases risk, and thus requires stronger evidence <sup>(39)</sup> [intuitionlabs.ai](#) <sup>(4)</sup> [emergingaihub.com](#)).

In practice, sponsors should document the question of interest early and include any supporting context or historical data. For example, if a model predicts patient dropout risk, one might point to prior trial data on dropout rates. The question then shapes the context of use, model design, and evaluation plan. As the Emerging AI Hub notes, framing the question precisely influences downstream planning: “the specificity matters because everything downstream flows from how you’ve framed this question” <sup>(2)</sup> [emergingaihub.com](#)). Without a clear question, validation metrics and risk analysis cannot be properly defined.

## Step 2: Define the Context of Use (COU)

Once the question is set, **Step 2** is to define the **Context of Use (COU)**: how exactly the AI model will operate to address that question <sup>(10)</sup> [emergingaihub.com](#) <sup>(41)</sup> [regulations.justia.com](#)). The COU adds the operational details around the reference question. It answers: What data will the model use? At what point in the process? Who sees the output and does what with it? For example, a COU description might specify: “The AI model will process ECG and lab data for each

screening patient, output a binary flag of 'eligible' or 'not eligible', and an investigator will be notified to make the inclusion decision."

The COU component is essential for scoping the evaluation. It includes:

- **Role of AI:** Is the model a decision-support tool or an automatic classifier? Will its output be reviewed by a human, or will it feed directly into an automated control system?
- **Boundaries and Inputs:** What inputs and training data are permissible? What patient population or batch conditions?
- **Integration with other evidence:** Is the AI the sole evidence for a decision, or one piece among many?

FDA guidance and commentators stress that understanding the AI's role (particularly the *degree of autonomy*) is crucial. For example, if an AI simply screens patients and a doctor makes the final call, the risk picture differs from an AI that autonomously enrolls or disqualifies patients. The Advisory Guidance clarifies that COU "includes the model's scope, its inputs/outputs, and its role relative to other evidence" (<sup>[42]</sup> intuitionlabs.ai). Accurately describing COU provides the baseline for all subsequent risk assessment and validation.

Several illustrative contexts of use are mentioned in the guidance, including:

- **Clinical Trial Design and Management** (e.g. optimizing inclusion criteria)
- **Patient Evaluation** (e.g. medical imaging or biomarker interpretation)
- **Endpoint Adjudication** (e.g. applying a digital biomarker to determine a clinical outcome)
- **Analyzing Trial Data** (e.g. meta-analytic or signal detection tools)
- **Pharmacovigilance** (e.g. AI triage of adverse event reports)
- **Manufacturing and Quality Control** (e.g. ensuring fill-volume or sterility using sensors) (<sup>[43]</sup> www.jdsupra.com).

The context description also flags uses outside the scope. The FDA guidance explicitly excludes flows such as early *discovery* (where molecules are generated by AI but then evaluated by traditional methods) and *operational* uses that do not directly affect safety or quality (e.g. scheduling or automated report writing) (<sup>[7]</sup> regulations.justia.com) (<sup>[8]</sup> emergingaihub.com). By focusing on COU, sponsors can confirm whether their AI use needs to follow these FDA recommendations (if it impacts regulatory decisions) or can use more proprietary development (e.g. trade secret models not submitted).

## Step 3: Assess Model Risk (Influence and Consequence)

With the COU in hand, **Step 3** is to **assess the model's risk level**, combining *model influence* and *decision consequence* (<sup>[4]</sup> emergingaihub.com) (<sup>[23]</sup> regulations.justia.com). This determines how stringent the rest of the credibility activities must be. The FDA proposes a two-dimensional risk matrix:

- **Model Influence:** how much the AI's output drives the decision. Is it one input among many (low influence), or the primary decision-maker (high influence)? For instance, an AI that generates a safety alert which a physician can override is lower influence; an AI that directly classifies a product batch pass/fail is high influence.
- **Decision Consequence:** how severe are the outcomes if the model is wrong? Does an error merely inconvenience the process, or does it risk patient harm or defective products? A misclassification of an insignificant lab value has low consequence; misidentifying a malignant tumor as benign (or failing to detect it) would be very high consequence.

The combination yields a risk rating (low, medium, high, etc.). A "high-risk" AI use might be high influence *and* high consequence. For example, the guidance gives a scenario where an AI assigns trial participants to home care vs. inpatient monitoring. If the model incorrectly deems a patient 'low risk' who is actually at high risk, the patient could suffer

serious harm. This is a high-influence, high-consequence context, demanding intense scrutiny (<sup>[4]</sup> [emergingaihub.com](#)). In contrast, an AI that, say, optimizes an administrative scheduling task (low consequence, low influence) would require minimal evidence.

Sponsors should **document their risk assessment rationale**. FDA notes that if the COU changes – e.g. the model is applied to a new patient population or the oversight differs – the risk analysis must be revisited (<sup>[4]</sup> [emergingaihub.com](#)). The Emerging AI Hub emphasizes that risk assessment is not a one-time checkbox but a **living document**: assumptions about influence and consequence should be recorded and updated as needed (<sup>[4]</sup> [emergingaihub.com](#)). In practice, this might take the form of a written justification or risk matrix table, explaining why each factor was scored as it was, so FDA reviewers can understand the sponsor's judgment.

Once risk is classified, it sets the **scope of credibility activities**. A high-risk model will trigger deeper investigations (e.g., extensive testing across edge cases, stricter performance criteria) and require submission of details such as code and raw data if feasible (<sup>[5]</sup> [www.jdsupra.com](#)) (<sup>[6]</sup> [pmc.ncbi.nlm.nih.gov](#)). A lower-risk model might need only summary results and supporting descriptions. But in all cases, sponsors must show that the level of evidence is commensurate with their risk categorization and COU.

► **Table: Key Documentation and Validation by Risk Level (Summary)**

## Development and Validation of AI Models in Drug Development

Once the question, COU, and risk are defined, sponsors advance to **validating the AI model** for its intended use (Steps 4–5). This involves rigorous testing of model inputs, algorithms, and outputs, and collecting evidence that the model performs reliably within its COU. Important domains of validation include:

### Data Quality and Fitness for Use

**High-quality training and test data** are the foundation of credible AI. The guidance and experts stress that model results are only as good as the data they're trained on. Sponsors must demonstrate that their data are **"fit-for-use"** in the given context (<sup>[44]</sup> [www.jdsupra.com](#)) (<sup>[9]</sup> [emergingaihub.com](#)). This includes:

- **Relevant and Representative Datasets.** The training set should reflect the COU. For instance, a model for patient selection should be trained on data from a similar patient population (similar demographics, disease severity, etc.). If important subpopulations are missing, the model may not generalize. Sponsors should describe data sources (e.g. clinical trial registries, electronic health records, manufacturing sensors) and ensure they cover the scope of use.
- **Data Preprocessing and Cleaning.** Any cleaning steps (handling missing values, standardizing measurements, removing outliers) should be documented. Explain how and why data were filtered or transformed. For example, if the model uses lab values, describe normalization methods.
- **Bias and Variance Checks.** Because historical data may contain biases, sponsors should check and document any potential biases in the training data. For example, if an ML model classifies patient risk, ensure that the dataset did not underrepresent a particular demographic. The Emerging AI Hub advises including bias and fairness assessment in the plan (<sup>[9]</sup> [emergingaihub.com](#)). Practical steps can include subgroup analysis (performance by age/gender/race) or techniques to de-bias data.
- **Data Integrity and Traceability.** Sponsors should maintain audit trails showing data provenance – where each dataset came from and how it was processed. This aligns with FDA's requirements for Good Documentation Practices. As a pharmaceutical manufacturing study notes, sustainable AI deployment "requires data monitoring systems that provide thorough tracking with complete traceability" (<sup>[25]</sup> [www.sciencedirect.com](#)). In regulated environments (GxP), each data transformation should be logged.

In regulatory submissions, sponsors will likely need to **describe the data** used to build the model. This could include data dictionaries, summaries of key variables, and sources. They should explain why these data are sufficient: e.g., “our training data encompassed 10,000 patients from four international trials, closely matching the Phase 3 study population”. Any limitations or assumptions (such as ignoring rare events) should also be acknowledged, as these affect credibility.

## Model Development and Architecture

Sponsors should detail how the AI model was built (Step 4). This covers the selection of algorithms, hyperparameters, model complexity, and software environment. While proprietary methods may remain confidential, the FDA guidance implies that *high-risk* models may require disclosure of the architecture. For example, the guidance notes disclosing model architecture and training methods “for high-risk AI models” (<sup>[5]</sup> [www.jdsupra.com](http://www.jdsupra.com)). Even if exact code is not submitted, sponsors should document:

- **Algorithmic Approach.** What type of model is used (e.g. random forest, neural network, gradient boosting)? Why was this approach chosen for the COU?
- **Feature Selection.** Which input variables feed into the model, and how were they selected? (In images, this might be pixel data or features; in tabular data, which lab tests, vitals, etc.).
- **Hyperparameters and Training Process.** General description of how the model was trained: scores used, regularization to avoid overfitting, random seeds for reproducibility, number of training epochs, etc. While not all details need public disclosure, internal documentation should record the training process.
- **Data Partitioning.** How was data split into training/validation/test sets? Was cross-validation used? For example, in trial modeling: “We held out 20% of data from trial A as a validation set, and tested the final model on independent Phase 2 data from trial B.” This helps FDA see that the model was not just tuned to one dataset.

If the AI is a combination or hybrid model (e.g. combining an empirical physical model with an ML component), sponsors should describe this hybrid structure. The Ark Biotech example illustrates a “hybrid” model for a bioreactor: the company reports that it integrated mechanistic reactor equations with ML to predict antibody yield (<sup>[17]</sup> [www.ark-biotech.com](http://www.ark-biotech.com)). In such cases, clarifying how the pieces interact is important, as the joint credibility depends on both parts.

Foley LLP notes that extensive transparency is needed: “comprehensive details regarding the AI model’s architecture, training methodologies, validation processes, and performance metrics” may have to be submitted (<sup>[5]</sup> [www.jdsupra.com](http://www.jdsupra.com)) for high-risk uses. In practice, sponsors should prepare technical documentation (like a model development report) that includes block diagrams or flowcharts of the model, even if not submitted (copies may be requested). Public summaries can be provided with as much detail as possible without revealing proprietary secrets.

## Performance Validation (Steps 4–5)

With the model in hand, sponsors must **execute the validation plan**. This means rigorously testing the model and quantifying its performance. The FDA expects a thorough demonstration that the model works as intended in the COU. Key elements include:

- **Independent Test Data.** Performance metrics should be computed on data not used for training or tuning. Ideally, this is a held-out dataset or separate study data. For example, [50] emphasizes using “truly independent test data”. If available, an external validation on data from a different source (another trial or site) greatly strengthens credibility.
- **Performance Metrics.** Choose metrics appropriate to the decision task. Common metrics include accuracy, sensitivity/specificity, positive predictive value, AUC-ROC, mean squared error, etc. For classification tasks, a confusion matrix is useful. For probabilistic outputs, calibration plots or Brier scores can show reliability. Whatever metrics are used, sponsors should justify why they align with the regulatory goal. In a safety context, one might prioritize sensitivity (catching all true positives), whereas in a resource-constrained context, one might balance specificity.

- **Statistical Uncertainty.** It is crucial to quantify uncertainty in performance estimates. Provide confidence intervals or statistical significance for metrics. For example, if an accuracy is reported as 92%, also state the 95% CI (e.g.  $\pm 2\%$ ). If possible, use bootstrapping or other resampling to estimate variability. This shows the range of possible performance and is part of “credibility evidence” (<sup>[1]</sup> [regulations.justia.com](https://www.regulations.gov)).
- **Bias and Fairness Analysis.** Demonstrate that performance is consistent across relevant subgroups. If the model will be applied to all adults, show that it works similarly for men and women, or across age groups, etc. Any systematic discrepancies should be investigated. The guidance underscores bias checks as part of risk management (<sup>[9]</sup> [emergingaihub.com](https://www.emergingaihub.com)). Note: absence of bias is not required for credibility, but sponsors must show they analyzed it.
- **Corner Cases and Robustness.** Test how the model performs in edge cases. For example, intentionally input unusual but plausible data to see if the model fails. Document any failure modes. Such robustness tests (sometimes called “sensitivity analyses”) can uncover weaknesses. Since AI models can sometimes give unpredictable results outside their training range, it’s prudent to test inputs marginally outside the expected distribution.
- **Comparison to Alternatives.** When relevant, compare the AI model’s performance to traditional methods. If a conventional algorithm or human decision is the gold standard, show that the AI is at least as good (or better). This contextualizes the benefit of AI. For example, if AI predicts clinical outcomes, compare against logistic regression or physician opinion.
- **Reproducibility Checks.** Run the model multiple times (if stochastic) to show results are stable. Save random seeds or apply techniques (e.g. fixed seeds or deterministic training pipelines) that ensure reproducibility of results, especially if code is audited.

All results should be meticulously **documented** (Step 5). As Emerging AI Hub notes, “every decision made during execution needs to be captured in real time. This is contemporaneous documentation, the kind that 21 CFR Part 11 and GxP environments already demand” (<sup>[45]</sup> [emergingaihub.com](https://www.emergingaihub.com)). In practice, this means storing audit logs of each model run, versions of code and data used, and signed conclusions. If electronic records are submitted, they must comply with Part 11 (audit trails, signatures) (<sup>[46]</sup> [emergingaihub.com](https://www.emergingaihub.com)).

Collectively, the validation evidence populates the credibility assessment report (Step 6). For example, Podichetty et al. suggest that high-risk AI submissions should include “model architecture, training logs, and data processing pipelines” in the documentation (<sup>[6]</sup> [pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov)). In simpler terms, the FDA expects to see “(1) the model, (2) the data used to develop the model, (3) model training, and (4) model evaluation including test data and performance metrics” (<sup>[47]</sup> [www.jdsupra.com](https://www.jdsupra.com)). These four components align with common regulatory sections (methods, data, validation plan, results) in an application, supplemented by AI-specific details.

## Cross-Cutting Considerations: Uncertainty, Bias, and Explainability

Beyond basic validation, the FDA framework encourages sponsors to address certain cross-cutting issues that bear on credibility:

- **Uncertainty Quantification.** As noted, any single metric is an estimate, and AI models have inherent uncertainty. Sponsors should propagate and report uncertainty. This might include confidence intervals on predictions or probability outputs. For instance, if predicting a lab value, provide not just a point prediction but also a prediction interval. Techniques like Bayesian modeling or ensemble methods can yield measures of confidence. Transparency about uncertainty is important for FDA reviewers to judge reliability.
- **Algorithmic Bias.** ML models can inadvertently perpetuate biases present in training data. The FDA explicitly mentions assessing bias as part of credibility planning (<sup>[9]</sup> [emergingaihub.com](https://www.emergingaihub.com)). For example, if an AI model for dermatology was trained mostly on light-skinned patients, its credibility would be suspect on darker skin. Sponsors should analyze performance by subgroup (race, gender, age, etc.) and either enrich the data or qualify the model’s limits. If biases are found, discuss mitigation strategies (reweighting, separate models, disclaimers). Documenting that the model does not have clinically significant biases (or explaining any discovered bias) strengthens credibility.

- **Explainability and Transparency.** While not a formal requirement in drug applications, explainability aids trust. The draft guidance does not mandate full interpretability for all AI, but Foley points out that one innovation need is “model explainability” (<sup>[14]</sup> [www.foley.com](http://www.foley.com)). Providing model rationales (feature importance, attention maps) can help FDA reviewers and investigators accept the results. Sponsors should at least document how the model makes decisions at a high level. For instance, decision-tree models have natural interpretability; neural networks might use SHAP values or similar.
- **Regulatory Compliance (Part 11, GxP).** Any software generating output for regulatory use is subject to FDA regulations for electronic records (21 CFR Part 11). This means maintaining audit trails for model runs, access controls, and validated software environments (<sup>[46]</sup> [emergingaihub.com](http://emergingaihub.com)). Sponsors should ensure AI tools are integrated into their quality systems. For example, if an AI is part of a Manufacturing Execution System (MES), it must comply with GMP record-keeping.

## Step 4–7: Planning, Execution, Documentation, and Review

Following risk assessment, Steps 4–7 formalize how to demonstrate credibility:

- **Step 4: Develop Credibility Plan.** Sponsors create a structured plan addressing all aspects of evidence collection. It should cover data quality checks, modeling strategy, performance criteria, validation design, and bias analysis (<sup>[9]</sup> [emergingaihub.com](http://emergingaihub.com)). Known as a “credibility assessment plan”, this document is akin to a validation protocol in clinical trials. It should specify acceptance criteria (e.g. “sensitivity  $\geq$  90% with 95% CI enclosed in trial standard”), data requirements, and responsibilities. The DLA Piper analysis emphasized that the entire framework’s purpose is to “establish trust in the AI model by collecting credibility evidence...for its particular COU” (<sup>[36]</sup> [www.dlapiper.com](http://www.dlapiper.com)), and the plan lays out exactly how that evidence will be obtained.
- **Step 5: Execute the Plan.** As discussed above, this is the implementation phase – running the model and analyses. All results are recorded. Especially for high-risk models, running the model in conditions mimicking the actual deployment (e.g. same software/hardware, same data pipelines) is recommended. Per 21 CFR 11, electronics records of these runs should be stored securely so FDA can audit them if needed (<sup>[46]</sup> [emergingaihub.com](http://emergingaihub.com)). The sponsor should track any deviations: e.g. if data quality issues forced more preprocessing than anticipated, or if a test dataset proved inadequate and a new one had to be obtained.
- **Step 6: Document Results and Deviations.** After testing, the sponsor compiles a **Credibility Assessment Report**. This report summarizes all planned analyses along with their outcomes and notes any deviations. The FDA expects “comprehensive documentation including results of all planned analyses, any deviations from the original plan, and the rationale” (<sup>[13]</sup> [emergingaihub.com](http://emergingaihub.com)). In other words, *nothing is hidden*. If a test failed to run or a metric fell short of the target, the sponsor must note it and discuss—preferably with a plan to address it. The document must allow an FDA reviewer to “independently evaluate the model’s credibility without needing to rerun any analyses” (<sup>[13]</sup> [emergingaihub.com](http://emergingaihub.com)). Thus, clear graphs, tables, and narrative explanation are essential. Providing code or extremely detailed appendices is optional, but clarity is mandatory.
- **Step 7: Determine Adequacy for COU.** Finally, the sponsor must judge whether the evidence shows the model is adequate. Adequacy is not a binary pass/fail, but a judgment call considering the model’s demonstrated performance in light of risk (<sup>[11]</sup> [emergingaihub.com](http://emergingaihub.com)). For example, a model with 95% accuracy might be adequate for a low-risk contingency plan but inadequate if missing the 5% leads to serious patient harm. The sponsor should explicitly state its conclusion: “Based on the evidence, the model is (or is not) sufficiently credible to support [the COU].” If it is not, additional work (e.g. more data, model revision) is needed before submission.

Importantly, Step 7 also brings **lifecycle considerations** into focus. An AI model’s performance may degrade over time if patient populations change or if process drifts. The guidance implies that sponsors should plan for **ongoing monitoring** of AI performance in post-approval or late-phase settings (<sup>[11]</sup> [emergingaihub.com](http://emergingaihub.com)) (<sup>[25]</sup> [www.sciencedirect.com](http://www.sciencedirect.com)). In effect, even after initial validation, the model may need re-evaluation if inputs shift. Although details of post-market AI plans are not submitted like trial data, sponsors should have standard operating procedures (SOPs) for periodic checking and re-validation. For example, if new assay equipment generates inputs differently, the model might need recalibration. The Emerging AI Hub notes that an “adequate model today may become inadequate if the data it processes shifts” (<sup>[11]</sup> [emergingaihub.com](http://emergingaihub.com)). Sponsors should therefore document a **change management plan** – how they will detect and respond to drift (e.g. thresholds for triggering re-training).

By completing these steps and assembling a thorough report, sponsors can demonstrate to FDA that they have “systematically built, documented, and evaluated” the model’s credibility (<sup>[48]</sup> [emergingaihub.com](http://emergingaihub.com)). In sum, establishing

credibility means **multiple layers of evidence**: strong data, robust performance testing, risk justification, and transparent reporting. The process mirrors that for any critical tool in drug development, with added attention to AI-specific issues.

## Data Analysis and Evidence – Building Credibility

Beyond the structured steps, it is useful to examine the underlying **data analysis practices** that strengthen an AI model's credibility. Here we consider how to handle data, measure performance, and interpret results in depth.

### Ensuring High-Quality Training Data

- **Data Source Evaluation:** Document the origin of every dataset used. Did it come from historical trials, real-world registries, manufacturing sensors, or other sources? Ensure these sources are reliable and relevant. For example, if training a model for a new diabetes drug trial, using EHR data from diabetic patients (rather than general population data) makes the model more credible.
- **Data Completeness and Labeling:** Check for missing or mislabeled data. For supervised models, labels must be accurate (e.g. true outcomes). If labels come from imperfect tests or human adjudicators, note their accuracy. Missingness should be handled systematically (e.g. imputation methods). The sponsor should report how many records were excluded for missing data or poor quality.
- **Preprocessing Procedures:** Any normalization, encoding (e.g. one-hot for categories), or filtering applied to the data must be described. For instance, continuous lab values might be log-transformed; imaging pixels might be standardized. Provide rationales (e.g. "log-transformation reduced skewness"). This assures FDA reviewers that the data fed into the AI were reasonably conditioned.
- **Feature Engineering:** If new features are derived (e.g. combining blood pressure and heart rate into a risk score), explain how. Does the model purely learn from raw inputs, or were domain transformations applied? Documenting feature engineering helps others understand model drivers.
- **Training/Test Splits:** Clearly specify how data were partitioned. A common best practice is 70% training, 30% testing, or k-fold cross-validation. Identify any external datasets if used for testing. For example, "Data from Trials A/B were used for training; Trial C data held out for final validation." This prevents "data leakage" where test data inadvertently influence training.
- **Size and Diversity:** State the sample size and diversity of the training data. AI models often require large datasets to generalize. If data are limited, consider data augmentation (for images) or note that limitations exist. For credible submission, extreme paucity of data should be justified and perhaps mitigated by conservative performance claims.

### Performance Evaluation

- **Metrics Selection:** Choose metrics that reflect the regulatory question. For categorical outcomes, measures like sensitivity (recall), specificity, and predictive values are typically reported. For regression tasks (e.g. predicting a concentration), use mean absolute error or  $R^2$ . Complex outcomes (like survival) might require concordance indices or Kaplan-Meier calibration plots. Always tie metrics to clinical meaning.
- **Summary Statistics:** Provide the mean, median, standard deviation, etc., of the model's outputs on test sets. For classification, provide a confusion matrix to show true positive/negative and false positive/negative counts. For probabilities, consider calibration plots. Clear tables of results allow others to gauge performance at a glance.
- **Visualizations:** Whenever possible, include figures such as ROC curves, precision-recall curves, or predicted-vs-actual scatter plots. These visuals communicate performance nuances and can highlight any systematic issues. For instance, an ROC curve can show trade-offs between sensitivity and specificity.
- **Uncertainty and Confidence:** As noted, give confidence intervals or p-values to accompany point estimates. FDA reviewers will appreciate knowing the statistical confidence of performance claims. If a subgroup analysis is done (e.g. model accuracy by age group), give separate intervals.

- **Benchmark Comparisons:** If relevant, compare the AI model against existing benchmarks or conventional methods. For example, show that a neural network outperforms logistic regression or a simpler rule-based algorithm. If AI is replacing human review, compare to inter-observer variability. Demonstrating superiority (or at least equivalence) to the status quo helps justify novelty and credibility.
- **Ad Hoc and Sensitivity Analyses:** Document any additional analyses done to test model stability. For example, "We retrained the model 100 times with different random seeds; accuracy varied by  $\pm 1\%$ . We also tested the model by adding 5% noise to inputs to simulate measurement error; performance only dropped from 92% to 89%." These sensitivity checks reveal robustness.

## Dealing with Bias and Fairness

Bias in AI is a major concern. Sponsors should proactively address it:

- **Identify Protected Attributes:** Determine which demographic or medical subgroups are relevant (race, sex, age, disease subtype, etc.).
- **Stratified Performance:** Evaluate and report model performance for each subgroup. For example, present accuracy for male vs. female, or for each ethnic group. If a model fares significantly worse in any subgroup, investigate why.
- **Balance or Reweight Data:** If a disparity is found due to underrepresentation (e.g. only 10% of training data are from group X), consider augmenting data or using techniques like SMOTE. Discuss any steps taken to balance.
- **Fairness Metrics:** Consider metrics like equalized odds or demographic parity if applicable. While not yet standardized in drug regulation, showing awareness of fairness principles (for example, by ensuring false positive rates are similar across groups) strengthens credibility.
- **Transparency of Limitations:** If no mitigation is possible (e.g. simply no data exist for a subgroup), clearly note it. FDA reviewers need to know the limits of generalization. In some cases, the sponsor may restrict the COU (e.g. not claiming applicability to under-tested groups).

## Documentation of Analysis

All analysis steps should be **logged and reproducible**:

- **Code and Environment:** Maintain the exact code (scripts, notebooks) used for training and testing. Although submitting proprietary code may not be required, having a repository with version control (e.g. Git) provides traceability. Record software versions of libraries and tools, so that an exact environment can be recreated (or simulated by FDA if needed).
- **Randomness Control:** If the model training has stochastic elements (random weight initialization, random data splits), ensure results are reproducible by fixing seeds or documenting the random state. This is part of *practical reproducibility*.
- **Record of Changes (Change Log):** Document every change made to the model after initial training. For example, if hyperparameters were tweaked, note the change and its effect. This log of modifications will feed into the credibility report as the project evolves.
- **Quality Management Integration:** Ideally, AI development is integrated into the sponsor's computer system validation and quality management (Part 11, GxP). For example, as the Emerging AI Hub notes, if AI outputs produce electronic records, Full 21 CFR 11 compliance (audit trails, signatures) is needed (<sup>[46]</sup> [emergingaihub.com](https://emergingaihub.com)). Ensuring that model training runs are captured in a validated environment underscores trust.

## Real-World Case Studies

Illustrative examples help ground these principles. The following cases—drawn from industry reports and regulatory submissions—show how AI credibility assessment can play out in practice.

### Case 1: AI in Bioprocess Optimization (Ark Biotech)

**Scenario:** A biopharmaceutical company (Ark Biotech) uses an AI-driven “hybrid” simulation model to optimize a monoclonal antibody production process. The model combines mechanistic bioreactor equations with machine learning elements to predict yield and quality based on process parameters (<sup>[17]</sup> [www.ark-biotech.com](http://www.ark-biotech.com)) (<sup>[18]</sup> [www.sciencedirect.com](http://www.sciencedirect.com)).

**COU and Question:** The model's COU is to predict whether a given set of critical process parameters (CPPs) in the upstream cell culture process will meet predefined Critical Quality Attribute (CQA) targets for the final antibody product. The regulatory question is essentially “Can we trust this model to replace some physical batch trials for CPP regime optimization?”

**Risk Assessment:** Because the model influences manufacturing decisions (e.g. adjusting nutrient feeds) and affects product quality (antibody potency, impurity levels), the risk is high. A wrong prediction could produce out-of-spec material. The company therefore applied the FDA's 7-step framework to ensure rigorous validation (<sup>[17]</sup> [www.ark-biotech.com](http://www.ark-biotech.com)).

**Data and Validation:** The sponsor gathered extensive historical bioreactor data and quality testing results from past batches as *credibility evidence*. These included measured CQAs for various CPP settings. They defined performance criteria (e.g. prediction accuracy within  $\pm 5\%$  of observed yields) and tested the model accordingly (<sup>[17]</sup> [www.ark-biotech.com](http://www.ark-biotech.com)). In practice, the team ran the model on a reserved hold-out set of batch data and computed prediction errors and confidence intervals, showing strong agreement. They also compared the hybrid AI predictions to a purely mechanistic model, demonstrating improved accuracy.

**Documentation:** Ark's report explicitly links the model output to a filed variation (CPP change) in their regulatory submission. They documented the COU (“predict product quality attributes under new CPPs”), the model's structure, the data sources, and their validation plan (<sup>[17]</sup> [www.ark-biotech.com](http://www.ark-biotech.com)). The credibility assessment plan included performance metrics (e.g. mean relative error) and bias checks (model performance across different cell batches). Any deviations (such as a subset of data with noisy measurements) were noted. By treating the model like a virtual assay, the sponsor effectively used AI to support a regulatory change, demonstrating the model's trustworthiness through documented evidence.

This example shows how an AI model in manufacturing can align with traditional quality-by-design practices. By embedding the AI validation into the pharma quality system (e.g. as part of process validation), the company addressed FDA expectations on documentation and traceability (<sup>[25]</sup> [www.sciencedirect.com](http://www.sciencedirect.com)). It also highlights that even clever “white box” or hybrid models require the same diligence as data-driven ones: the path from inputs to prediction must be validated and transparent.

## Case 2: AI for Patient Selection (COVID-19 Trial)

**Scenario:** Early in the COVID-19 pandemic, an AI research team trained a machine-learning model to identify severe inflammatory signatures in hospitalized patients, predicting who would benefit from the immunomodulatory drug anakinra. The model analyzed combinations of lab values and clinical data to stratify patients by likely response (<sup>[19]</sup> [pmc.ncbi.nlm.nih.gov](http://pmc.ncbi.nlm.nih.gov)).

**COU and Question:** The model's COU was patient stratification under an FDA Emergency Use Authorization (EUA) context. The question was: “Which COVID-19 patients at hospital admission are most likely to benefit from anakinra treatment in terms of survival or disease progression?” This directly informed treatment decisions, making it a high-consequence use.

**Risk Assessment:** Patient selection is high risk: a false negative (missing a patient who would benefit) could worsen outcomes, while a false positive could expose patients to unnecessary treatment. Therefore, the model had to be very reliable. FDA EUA submissions for COVID therapeutics required strong evidence, and adding AI meant meeting analogous standards to biomarker qualification.

**Data and Validation:** The model was trained on retrospective patient data (lab tests, demographics, outcomes) from hospitals in Italy. As publicized in *Clinical Pharmacology & Therapeutics*, the team reported that the AI model correctly identified a subset of patients who showed markedly improved outcomes with anakinra under treatment<sup>[19]</sup> ([pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov)). Performance was validated on separate patient datasets and showed statistically significant predictive power. The model's positive predictive value was high, and confidence intervals were calculated.

**Documentation:** In the EUA submission, the evidence would have included a description of the model (features used, algorithm type), training data characteristics, and validation results. The credibility narrative emphasized targeted patient selection ("supporting targeted treatment strategies" (<sup>[19]</sup> [pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov))). Although FDA internal details are scarce, the published account demonstrates how the model effectively became part of the clinical trial design. The team likely provided FDA with the model's decision criteria (essentially a calculated risk score) and accompanying data analysis showing benefit in trials.

This real-world case shows AI assisting in **adaptive trial design** and precision medicine. It illustrates how an AI model's predictions were vetted and used in an actual regulatory context. It also underscores the need for clarity: the published tables from FDA (NIH Clinical Trials registry) and articles served as de facto documentation of the AI's efficacy. The model essentially functioned as a digital biomarker. If this submission had been part of an NDA or EUA, a credibility report would review all this analysis step by step.

### Case 3: AI-Enabled Image Analysis for Endpoints (NASH Trial)

**Scenario:** A biotech company developed an AI algorithm for histopathology: it quantifies liver biopsy features in nonalcoholic steatohepatitis (NASH) trials. The model analyzes digitized histology images to produce a quantitative score of disease activity. This was submitted as a "Qualified Biomarker" by the European Medicines Agency (EMA) and under discussion with the FDA.

**COU and Question:** The COU is to provide an *objective endpoint* in clinical trials: "Is this patient's liver improving?" The question was, effectively, "Can the AI reliably score disease severity and detect changes over time in NASH trials?" The AI was used alongside or instead of pathologist readings to measure treatment effect.

**Risk Assessment:** The risk is moderate-to-high. The endpoint of a clinical trial is critical for proving efficacy. However, human histology scoring itself is known to have inter-reader variability, so a consistent AI could reduce noise. The main concern was ensuring the AI's quantitative measurements are reproducible and not biased by artifacts (e.g. scanner differences).

**Data and Validation:** The developers trained the model on thousands of biopsy images labeled by expert pathologists. Validation involved an independent set of biopsies, where the AI's scores were compared against a panel of pathologists. The model showed high concordance and lowered variability. They also tested the AI on images from different labs to ensure it generalized.

**Documentation:** In the EMA Qualification Opinion, regulators reviewed the evidence of the model's performance. The opinion notes that the AI provided "reproducible, quantitative measurements that improve the consistency of disease activity assessment" (<sup>[20]</sup> [pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov)). The documentation included a description of the neural network architecture, training datasets, and validation results. It likely also addressed image preprocessing steps (color normalization, artifact removal) and robustness to slide variability. Because the AI underpinned a key trial endpoint, the documentation had to be especially rigorous—essentially matching the standard for a medical device.

This case illustrates a regulatory acceptance of an AI as part of the evidence package. It highlights that machine-learning tools can become *qualified biomarkers*, undergoing a credibility assessment akin to ligands or imaging assays. The thorough validation of the AI (including generalizability across image sources) and transparency about its performance enabled regulators to trust the output. For sponsors planning similar image-based AI, this sets a precedent: high-quality evidence and alignment with existing biomarker qualification frameworks are essential.

# Documentation and Regulatory Submission

A central theme of the FDA guidance is that all aspects of AI development must be **documented** for regulatory review. In a submission (e.g. IND, NDA, BLA filings), sponsors should include an AI addendum or sections in relevant places (clinical protocols, CMC sections, etc.) describing the AI model. Key documentation items include:

- **Model Description:** A clear narrative of what the model is and does (COU, question it answers). This may go in the clinical or quality sections depending on use. Provide a summary of its design and purpose.
- **Inputs and Outputs:** Define all inputs (features, sensors, etc.) and outputs. For instance, if the AI works on MRI scans, describe the imaging protocol and preprocessing. If it ingests blood test results, list them. Also specify how outputs are interpreted (e.g. risk score thresholds).
- **Training and Validation Data:** Describe the datasets used (size, source, features) for training and validation. If using proprietary datasets (e.g. from prior trials), explain their provenance and relevance. Indicate which data were internal vs. any external or public sources.
- **Performance Metrics and Results:** Provide tables/figures of model performance on test sets, including metrics (accuracy, sensitivity, etc.) and statistical confidence. Include subgroup analyses (e.g. by relevant patient characteristics). This helps demonstrate reliability.
- **Risk Analysis:** Summarize the risk assessment (Step 3). This can be a simple table or paragraph explaining the chosen influence/consequence levels and justification. It should show that the sponsor has considered the implications of model failure.
- **Validation Procedures:** Outline the specific validation studies performed. E.g. "The model was evaluated on a blinded external dataset of 200 patients, achieving X% accuracy. We also retrospectively applied the model to historical trial participants and confirmed that treatment effect was predicted correctly." Even if not full clinical trial validity, showing the model reproduces known findings boosts confidence.
- **Change Control Plan:** Describe how model updates will be handled post-submission. For example, FDA suggests manufacturers plan for "life cycle maintenance of credibility" (<sup>[49]</sup> [www.foley.com](http://www.foley.com)). This means having SOPs for model revisions: when new data accrue, the model may be optionally retrained, but in that case a regulatory amendment or inspection will cover it. Outline this plan in writing.
- **Software Validation:** If submitting software (e.g. algorithm as part of a device or analytical package), include evidence of software testing and qualification. Otherwise, at least assert that the computational infrastructure is validated for its intended use.
- **Transparency vs. IP:** As Foley and others warn, sponsors must balance the FDA's disclosure needs with intellectual property. The guidance implies that high-risk models will not easily remain secret (<sup>[15]</sup> [www.foley.com](http://www.foley.com)). Sponsors should consult legal counsel: often the recommendation is to patent innovative AI methods before submission to protect IP, since FDA will want to see training methods and code descriptions (<sup>[15]</sup> [www.foley.com](http://www.foley.com)) (<sup>[14]</sup> [www.foley.com](http://www.foley.com)). Where possible, use nondisclosure agreements or summary descriptions for especially sensitive details, acknowledging to FDA (and noting in submissions) any trade secret claims, but be prepared to give as much information as needed to establish credibility. For example, one might describe the model architecture in a patent application filed alongside the submission.
- **Integration with Quality Systems:** Document how the AI development project fits into the company's quality management. This might include references to GLP/GMP standards, computer system validation records, or good machine learning practice (GMLP) guidelines. Showing that the AI code and data are managed under QA controls earns trust.

Overall, the goal is for a regulatory reviewer to **follow the train of logic**: from COU to risk to validation to conclusion. The documentation package should include a written trace of all steps and decisions (much like a clinical study report). Reviewers should not have to guess how the model was tested or where results came from. As Emerging AI Hub emphasizes, transparency is key: "The documentation package should be structured so that an FDA reviewer can independently evaluate the model's credibility without needing to rerun any analyses" (<sup>[13]</sup> [emergingaihub.com](http://emergingaihub.com)).

## Implications, Challenges, and Future Directions

The FDA's AI credibility guidance represents a major milestone. It signals that AI/ML is maturing from experimental to regulated territory in drug development. However, it also raises significant implications and challenges:

- **Innovation Opportunities.** By spelling out requirements, the guidance creates clear targets for innovation. Foley identifies opportunities such as explainable AI techniques, automated bias-detection tools, and systems for model monitoring and reporting (<sup>[21]</sup> [www.foley.com](http://www.foley.com)). For example, "automatic systems to generate reports of model development and evaluation" are specifically cited, since sponsors will need to compile lots of documentation (<sup>[21]</sup> [www.foley.com](http://www.foley.com)). Companies that develop tools for AI transparency or life-cycle management (e.g. drift detectors, automated audit logs) could become valuable partners. Cryptographic methods (federated learning, homomorphic encryption) might help protect data IP under these disclosure rules.
- **Intellectual Property Tensions.** As discussed, the requirement for transparency may push firms to pursue patents rather than trade secrets for their AI innovations (<sup>[15]</sup> [www.foley.com](http://www.foley.com)) (<sup>[14]</sup> [www.foley.com](http://www.foley.com)). This is a strategic shift: AI techniques that once were kept proprietary will likely end up published. Early patent filing and careful redaction of genuinely non-public algorithms (if justifiable) will be critical. We expect sponsors will engage FDA early (pre-IND meetings, Type C meetings) to clarify what can remain confidential and how to align IP strategy with regulatory disclosure.
- **Regulatory Harmonization.** The FDA is not alone in addressing AI. The European Medicines Agency (EMA) published a reflection paper in September 2024 on AI in the entire medicine lifecycle ([www.ema.europa.eu](http://www.ema.europa.eu)), emphasizing a risk-based, human-centric approach. EmergingAIHub notes that in January 2026 FDA and EMA released joint principles for "good AI practice" in drug development (<sup>[22]</sup> [emergingaihub.com](http://emergingaihub.com)). There is also work from standards bodies (e.g. ISPE's July 2025 GAMP guide on AI in GxP). While details differ, a convergence is apparent: risk-based validation, transparency, and post-market monitoring are common themes. Sponsors should breathe a sigh of relief that global regulatory environments seem to be coalescing rather than diverging. Preparing for FDA's framework will likely put a company in good position for EMA or other regulatory bodies.
- **Operational Burden.** The compliance cost will be non-trivial. Building a comprehensive justification and documentation for an AI model rivals that for a drug formulation. Smaller biotech firms may find the burden heavy; they may need to partner with AI service providers who specialize in regulatory-grade ML. It will also require cross-functional teams: data scientists, statisticians, clinicians, and regulatory affairs experts must collaborate. Training programs on "regulatory ML" may emerge.
- **Data Considerations.** The emphasis on data quality highlights an ongoing need for better data infrastructure. Poor data can sink AI credibility. Regulators may therefore push for centralized data standards (like common formats for electronic health records) to ease future AI validation. Companies may invest more in data curation and bias audits, even outside regulated contexts.
- **Future FDA Actions.** This draft guidance will be open for comment, and the final version (likely in late 2025 or 2026) may refine details. FDA has already indicated interest in further guidance on AI in postmarketing safety. Stakeholders should watch for Class II combinations of AI-enabled devices with drugs, other emerging topics like generative AI in protocols, and guidance on specific subareas (e.g. pharmacovigilance AI).
- **Case Study Evidence Base.** As more AI applications reach submissions, FDA reviewers will accumulate experience. For example, the Emerging AI Hub reports over 1,060 submissions involving AI/ML through 2025 (<sup>[31]</sup> [emergingaihub.com](http://emergingaihub.com)) – a number that will continue growing. This means FDA staff are learning what works and what problems arise. The feedback cycle (comments on the draft) will also shape future policies. Sponsors should monitor FDA's published decisions (e.g. summary reviews) to glean how the credibility framework is being applied in reality.
- **Ethical and Societal Implications.** Beyond technicalities, credible AI in healthcare raises questions of trust and consent. Patients and clinicians may be wary of "black box" models. The emphasis on documentation and explainability can help build public trust. Moreover, transparent validation reports could end up as part of regulatory records accessible to healthcare professionals, not just regulators. In that sense, the guidance nudges the industry toward a more open model of algorithmic accountability.

## Conclusion

The FDA's draft guidance on AI credibility marks a pivotal evolution in drug regulatory science. It signals that AI models – once treated loosely as computational aids – must now meet rigorous evidence standards when they impact patient or product outcomes (<sup>[30]</sup> [www.fda.gov](http://www.fda.gov)) (<sup>[1]</sup> [regulations.justia.com](http://regulations.justia.com)). The seven-step framework provides a clear (if demanding) roadmap: define a specific question and context, analyze risks, and then validate the model to a degree commensurate with those risks (<sup>[1]</sup> [regulations.justia.com](http://regulations.justia.com)) (<sup>[2]</sup> [emergingaihub.com](http://emergingaihub.com)). Throughout, the focus is on *credible evidence*: high-quality data, robust testing, and meticulous documentation.

This report has unpacked those requirements in detail. We traced the historical impetus (ever-increasing AI use in submissions (<sup>[30]</sup> [www.fda.gov](http://www.fda.gov)) (<sup>[31]</sup> [emergingaihub.com](http://emergingaihub.com))), explained each step of the new guidance, and highlighted key data and validation practices. Multiple viewpoints (legal, technical, academic) and real examples demonstrate how to apply these principles in practice. Importantly, we saw that credibility is holistic: it is built not only on one metric but on a chain of rigor – from how you frame the question to how you monitor model drift after approval (<sup>[11]</sup> [emergingaihub.com](http://emergingaihub.com)) (<sup>[25]</sup> [www.sciencedirect.com](http://www.sciencedirect.com)).

In summary, evaluating and documenting an AI model for drug development means treating it like any critical scientific instrument. Sponsors must plan from the outset: choosing appropriate data, selecting performance goals, and recording every step. They must be ready to explain the AI to regulators comprehensively, often in ways they never had to for traditional statistical models. As FDA's Commissioner Califf noted, the goal is to foster innovation while upholding safety and efficacy standards (<sup>[28]</sup> [www.fda.gov](http://www.fda.gov)). By following the FDA's framework, developers can embrace AI's promise – using it to accelerate drug development and improve care – with the safeguards needed for public trust.

**Sources:** This report draws on FDA statements and draft guidance (<sup>[30]</sup> [www.fda.gov](http://www.fda.gov)) (<sup>[1]</sup> [regulations.justia.com](http://regulations.justia.com)), peer-reviewed analyses (<sup>[29]</sup> [academic.oup.com](http://academic.oup.com)) (<sup>[6]</sup> [pmc.ncbi.nlm.nih.gov](http://pmc.ncbi.nlm.nih.gov)), legal and industry commentary (<sup>[5]</sup> [www.jdsupra.com](http://www.jdsupra.com)) (<sup>[36]</sup> [www.dlapiper.com](http://www.dlapiper.com)), and relevant case studies (<sup>[17]</sup> [www.ark-biotech.com](http://www.ark-biotech.com)) (<sup>[19]</sup> [pmc.ncbi.nlm.nih.gov](http://pmc.ncbi.nlm.nih.gov)). All claims above are supported by cited evidence. Any quotations from FDA or expert sources are explicitly referenced with hyperlinks to the source.

---

## External Sources

- [1] <https://regulations.justia.com/regulations/fedreg/2025/01/07/2024-31542.html#:~:Speci...>
- [2] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:Step%...>
- [3] <https://www.jdsupra.com/legalnews/ai-drug-development-fda-releases-draft-3721117/#:~:The%2...>
- [4] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:Step%...>
- [5] <https://www.jdsupra.com/legalnews/ai-drug-development-fda-releases-draft-3721117/#:~:For%2...>
- [6] <https://pmc.ncbi.nlm.nih.gov/articles/PMC12675821/#:~:For%2...>
- [7] <https://regulations.justia.com/regulations/fedreg/2025/01/07/2024-31542.html#:~:model...>
- [8] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:What%...>
- [9] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:Step%...>
- [10] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:Step%...>
- [11] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:Step%...>
- [12] <https://www.sciencedirect.com/science/article/pii/S0022354926000894#:~:chrom...>
- [13] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:Step%...>
- [14] <https://www.foley.com/de/insights/publications/2025/01/ai-drug-development-fda-releases-draft-guidance/#:~:life...>
- [15] <https://www.foley.com/de/insights/publications/2025/01/ai-drug-development-fda-releases-draft-guidance/#:~:Paten...>
- [16] <https://www.foley.com/de/insights/publications/2025/01/ai-drug-development-fda-releases-draft-guidance/#:~:drift...>
- [17] <https://www.ark-biotech.com/insights/making-pharma-ai-ready-applying-the-fdas-draft-guidance/#:~:Be cau...>
- [18] <https://www.sciencedirect.com/science/article/pii/S0022354926000894#:~:The%2...>

- [19] <https://pmc.ncbi.nlm.nih.gov/articles/PMC12675821/#:~:Patie...>
- [20] <https://pmc.ncbi.nlm.nih.gov/articles/PMC12675821/#:~:reduc...>
- [21] <https://www.foley.com/de/insights/publications/2025/01/ai-drug-development-fda-releases-draft-guidance/#:~:subm...>
- [22] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:The%2...>
- [23] <https://regulations.justia.com/regulations/fedreg/2025/01/07/2024-31542.html#:~:frame...>
- [24] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:Step%...>
- [25] <https://www.sciencedirect.com/science/article/pii/S0022354926000894#:~:assis...>
- [26] <https://intuitionlabs.ai/articles/fda-ai-drug-development-guidance/#:~:Artif...>
- [27] <https://pmc.ncbi.nlm.nih.gov/articles/PMC12675821/#:~:covar...>
- [28] <https://www.fda.gov/news-events/press-announcements/fda-proposes-framework-advance-credibility-ai-models-used-drug-and-biological-product-submissions#:~:~%E2%8...>
- [29] <https://academic.oup.com/jlb/article-abstract/doi/10.1093/jlb/lsaf028/8316994#:~:The%2...>
- [30] <https://www.fda.gov/news-events/press-announcements/fda-proposes-framework-advance-credibility-ai-models-used-drug-and-biological-product-submissions#:~:~The%2...>
- [31] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:This%...>
- [32] <https://www.fda.gov/news-events/press-announcements/fda-proposes-framework-advance-credibility-ai-models-used-drug-and-biological-product-submissions#:~:~ln%20...>
- [33] <https://www.fda.gov/news-events/press-announcements/fda-proposes-framework-advance-credibility-ai-models-used-drug-and-biological-product-submissions#:~:~Today...>
- [34] <https://intuitionlabs.ai/articles/fda-ai-drug-development-guidance/#:~:~ln%20...>
- [35] <https://regulations.justia.com/regulations/fedreg/2025/01/07/2024-31542.html#:~:~Speci...>
- [36] <https://www.dlapiper.com/en-it/insights/publications/2025/01/fda-releases-draft-guidance-on-use-of-ai/#:~:~FDA%2...>
- [37] <https://regulations.justia.com/regulations/fedreg/2025/01/07/2024-31542.html#:~:~submi...>
- [38] <https://www.foley.com/de/insights/publications/2025/01/ai-drug-development-fda-releases-draft-guidance/#:~:~promo...>
- [39] <https://intuitionlabs.ai/articles/fda-ai-drug-development-guidance/#:~:~For%2...>
- [40] <https://intuitionlabs.ai/articles/fda-ai-drug-development-guidance/#:~:~case%...>
- [41] <https://regulations.justia.com/regulations/fedreg/2025/01/07/2024-31542.html#:~:~credi...>
- [42] <https://intuitionlabs.ai/articles/fda-ai-drug-development-guidance/#:~:~risk%...>
- [43] <https://www.jdsupra.com/legalnews/ai-drug-development-fda-releases-draft-3721117/#:~:~1,Lif...>
- [44] <https://www.jdsupra.com/legalnews/ai-drug-development-fda-releases-draft-3721117/#:~:~The%2...>
- [45] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:~This%...>
- [46] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:~:,syst...>
- [47] <https://www.jdsupra.com/legalnews/ai-drug-development-fda-releases-draft-3721117/#:~:~High,...>
- [48] <https://emergingaihub.com/ai-compliance/fda-ai-credibility-framework-guide/#:~:~frami...>
- [49] <https://www.foley.com/de/insights/publications/2025/01/ai-drug-development-fda-releases-draft-guidance/#:~:~ln%20...>

## IntuitionLabs - Industry Leadership & Services

**North America's #1 AI Software Development Firm for Pharmaceutical & Biotech:** IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**AI Chatbot Development:** Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**AI Consulting & Training:** Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at <https://intuitionlabs.ai/contact> for a consultation.

---

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by [Adrien Laurent](#), a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.

© 2025 IntuitionLabs.ai. All rights reserved.