# AI Tester Interview Questions & Structure

## Interview Structure (90-120 minutes)

1. **Introduction & Background** (10 minutes)
2. **Fundamental AI Concepts** (20 minutes)
3. **Testing Methodologies** (25 minutes)
4. **Technical Deep Dive** (25 minutes)
5. **Scenario-Based Questions** (20 minutes)
6. **Ethics & Bias** (15 minutes)
7. **Questions from Candidate** (5 minutes)

---

## Category 1: Fundamental AI Concepts

### Question 1: What is the difference between supervised and unsupervised learning?

| Classic AI | GenAI |
|---|---|
| **Answer:** Supervised learning uses labeled training data to learn input-output mappings (e.g., classification, regression). Unsupervised learning finds patterns in unlabeled data (e.g., clustering, dimensionality reduction). | **Answer:** In GenAI, supervised learning is used for tasks like fine-tuning models on specific datasets (e.g., instruction-following). Unsupervised learning is fundamental to foundation models trained on vast unlabeled text/image data to learn representations. |

### Question 2: How would you explain overfitting and how to detect it?

| Classic AI | GenAI |
|---|---|
| **Answer:** Overfitting occurs when a model performs well on training data but poorly on new data. Detection methods include validation curves, cross-validation, and monitoring training vs. validation loss. Solutions include regularization, dropout, and more training data. | **Answer:** In GenAI, overfitting manifests as models memorizing training data rather than generalizing. Detection includes perplexity metrics, human evaluation, and checking for verbatim reproduction of training examples. Solutions include diverse training data, regularization techniques, and careful fine-tuning. |

### Question 3: What are the key differences between traditional ML and deep learning?

| Classic AI | GenAI |
|---|---|
| **Answer:** Traditional ML uses handcrafted features with simpler algorithms (SVM, decision trees). Deep learning automatically learns hierarchical features through neural networks with multiple layers. Deep learning requires more data and computational power but can handle complex patterns. | **Answer:** GenAI builds on deep learning principles but focuses on generative capabilities. It uses transformer architectures, attention mechanisms, and massive scale. Unlike traditional ML's predictive focus, GenAI creates new content and handles multimodal tasks. |

## Category 2: Testing Methodologies

### Question 4: How would you design a test suite for an AI model?

| Classic AI | GenAI |
|---|---|
| **Answer:** Include unit tests for data preprocessing, integration tests for model pipeline, performance tests (accuracy, precision, recall), robustness tests (adversarial examples), and regression tests for model updates. Use train/validation/test splits and cross-validation. | **Answer:** Test suite includes functional tests (output quality, coherence), safety tests (harmful content detection), robustness tests (prompt injection, jailbreaking), performance tests (latency, throughput), and bias evaluation. Include human evaluation and automated metrics like BLEU, ROUGE. |

### Question 5: What metrics would you use to evaluate model performance?

| Classic AI | GenAI |
|---|---|
| **Answer:** Classification: accuracy, precision, recall, F1-score, AUC-ROC. Regression: MSE, RMSE, MAE, $R^2$. Also consider confusion matrices, learning curves, and domain-specific metrics. | **Answer:** Text generation: BLEU, ROUGE, perplexity, human evaluation scores. Image generation: FID, IS, CLIP score. General: coherence, relevance, factual accuracy, safety scores. Include both automated metrics and human judgment. |

### Question 6: How would you test for edge cases and corner cases?

| Classic AI | GenAI |
|---|---|
| **Answer:** Create boundary value tests, null/empty inputs, outlier data, adversarial examples, and stress tests with extreme values. Use synthetic data generation and mutation testing. Test with data outside training distribution. | **Answer:** Test with unusual prompts, multilingual inputs, edge case scenarios, prompt injection attempts, and boundary conditions. Include adversarial prompting, context length limits, and uncommon use cases. Test model behavior with conflicting instructions. |

# Category 3: Technical Deep Dive

## Question 7: Explain the concept of attention mechanisms.

| Classic AI | GenAI |
|---|---|
| **Answer:** Attention allows models to focus on relevant parts of input data, originally developed for sequence-to-sequence tasks. It creates weighted connections between input and output elements, solving the bottleneck problem of fixed-size representations. | **Answer:** Self-attention in transformers is fundamental to GenAI. It allows models to relate different positions in a sequence, enabling better understanding of context and long-range dependencies. Multi-head attention provides multiple representation subspaces. |

## Question 8: What is transfer learning and when would you use it?

| Classic AI | GenAI |
|---|---|
| **Answer:** Transfer learning uses pre-trained models as starting points for new tasks, especially useful with limited data. Common approaches include feature extraction and fine-tuning. Effective when source and target domains are related. | **Answer:** Transfer learning is core to GenAI through foundation models. Pre-trained models like GPT, BERT are fine-tuned for specific tasks. Includes techniques like few-shot learning, prompt engineering, and parameter-efficient fine-tuning (LoRA, adapters). |

## Question 9: How would you handle data quality issues in AI systems?

| Classic AI | GenAI |
|---|---|
| **Answer:** Implement data validation pipelines, outlier detection, duplicate removal, and consistency checks. Use data profiling, automated quality monitoring, and establish data governance processes. Handle missing values and feature engineering. | **Answer:** Focus on training data quality, content filtering, deduplication at scale, and bias detection. Implement safety filters, toxicity detection, and factual accuracy checks. Monitor for data leakage and privacy concerns in large datasets. |

# Category 4: Scenario-Based Questions

## Question 10: A model performs well in testing but fails in production. How would you investigate?

| Classic AI | GenAI |
|---|---|
| **Answer:** Check for data drift, feature distribution changes, infrastructure issues, or training/production environment mismatches. Implement monitoring, logging, and gradual rollout strategies. Analyze prediction confidence and error patterns. | **Answer:** Investigate prompt distribution changes, context length issues, safety filter interactions, or latency problems. Check for jailbreaking attempts, unexpected user patterns, or content policy violations. Monitor generation quality and user satisfaction. |

## Question 11: How would you test an AI system for bias?

| Classic AI | GenAI |
|---|---|
| **Answer:** Use fairness metrics (demographic parity, equalized odds), analyze performance across different demographic groups, check for historical bias in training data, and implement bias detection tools. Test with synthetic diverse datasets. | **Answer:** Test outputs across different demographic groups, cultural contexts, and languages. Check for stereotypical representations, unfair treatment of protected groups, and cultural biases. Use automated bias detection tools and human evaluation from diverse perspectives. |

## Question 12: Describe how you would implement A/B testing for AI models.

| Classic AI | GenAI |
|---|---|
| **Answer:** Split traffic between model versions, define success metrics, ensure statistical significance, and control for confounding variables. Monitor performance metrics, user behavior, and business outcomes. Implement gradual rollout and rollback capabilities. | **Answer:** Similar approach but consider generation quality, user engagement, safety metrics, and computational costs. Test different model sizes, prompting strategies, or fine-tuning approaches. Monitor for harmful outputs and user satisfaction in addition to performance metrics. |

# Category 5: Ethics & Bias

## Question 13: What are the main ethical concerns in AI testing?

| Classic AI | GenAI |
|---|---|
| **Answer:** Bias in training data and algorithms, privacy concerns, transparency and explainability, fairness across different groups, and accountability for decisions. Need for diverse testing teams and inclusive design processes. | **Answer:** Additional concerns include misinformation generation, deepfakes, copyright infringement, job displacement, and potential for harmful content creation. Need for content filtering, safety alignment, and responsible deployment practices. |

## Question 14: How would you ensure AI safety in your testing process?

| Classic AI | GenAI |
|---|---|
| **Answer:** Implement robust validation processes, test for adversarial attacks, ensure human oversight, establish clear failure modes, and maintain audit trails. Regular security assessments and compliance with regulations. | **Answer:** Include red team testing, safety alignment evaluation, harmful content detection, and robustness against prompt injection. Implement constitutional AI principles, human feedback integration, and continuous monitoring for emerging safety issues. |

# Extended Technical Questions

## Question 15: Explain different types of RAG architectures and their testing considerations.

| Classic AI | GenAI |
|---|---|
| **Answer:** N/A - RAG is specific to GenAI. For classic AI, similar concept would be information retrieval systems combined with traditional ML models, tested for retrieval precision/recall and model accuracy separately. | **Answer: Naive RAG**: Simple retrieve-then-generate. **Advanced RAG**: Pre-retrieval (query rewriting), retrieval (hybrid search), post-retrieval (re-ranking). **Modular RAG**: Flexible components. **Testing**: Evaluate retrieval quality (relevance, coverage), generation quality (faithfulness, coherence), and end-to-end performance. |

## Question 16: What is chunking in RAG and how would you test different chunking strategies?

| Classic AI | GenAI |
|---|---|
| **Answer:** N/A - Chunking is RAG-specific. In classic AI, similar concept is feature segmentation or data preprocessing, tested for information preservation and model performance impact. | **Answer: Chunking**: Splitting documents into manageable pieces. **Strategies**: Fixed-size, sentence-based, semantic, recursive. **Chunk overlap**: Prevents information loss at boundaries. **Testing**: Evaluate retrieval accuracy, information completeness, context preservation. Test different chunk sizes (256, 512, 1024 tokens) and overlap percentages (10-20%). |

## Question 17: Compare different fine-tuning approaches and their testing requirements.

| Classic AI | GenAI |
|---|---|
| **Answer:** Transfer learning approaches: feature extraction (freeze base layers), fine-tuning (adjust all layers), progressive unfreezing. Test for catastrophic forgetting, domain adaptation effectiveness, and computational efficiency. | **Answer: Full fine-tuning**: Update all parameters. **Parameter-efficient**: LoRA, adapters, prefix tuning. **Instruction tuning**: Task-specific formatting. **RLHF**: Human preference alignment. **Testing**: Evaluate task performance, general capability retention, computational cost, and training stability. |

## Question 18: Explain the pre-training process and how you would validate it.

| Classic AI | GenAI |
|---|---|
| **Answer:** Pre-training involves unsupervised feature learning on large datasets (e.g., autoencoders, word2vec). Validation includes representation quality, downstream task performance, and computational efficiency. | **Answer: Pre-training**: Large-scale unsupervised learning on diverse text/multimodal data. **Validation**: Perplexity metrics, downstream task evaluation, scaling laws verification, data quality assessment, training stability monitoring, and emergent capability evaluation. |

## Question 19: What are the key hyperparameters and how would you test their impact?

| Classic AI | GenAI |
|---|---|
| **Answer: Learning rate**, regularization strength, batch size, network architecture, optimizer choice. **Testing**: Grid search, random search, Bayesian optimization. Monitor training curves, validation performance, and generalization gaps. | **Answer: Learning rate**, batch size, sequence length, **temperature** (randomness in generation), top-p/top-k sampling, attention heads, model depth. **Testing**: Systematic hyperparameter sweeps, generation quality assessment, and computational cost analysis. |

## Question 20: What is temperature in language models and how would you test its effects?

| Classic AI | GenAI |
|---|---|
| **Answer:** N/A - Temperature is specific to generative models. In classic AI, similar concept is confidence thresholding or prediction uncertainty, tested for calibration and decision boundaries. | **Answer: Temperature**: Controls randomness in token selection. Low (0.1-0.3): deterministic, focused. High (0.8-1.0): creative, diverse. **Testing**: Evaluate output diversity, coherence, factual accuracy, and user preference across temperature ranges. Test consistency and controllability. |

## Question 21: Describe different neural network architectures and their testing considerations.

| Classic AI | GenAI |
|------------|-------|
| **Answer: CNNs**: Convolution, pooling, fully connected layers. **RNNs**: LSTM, GRU for sequences. **Testing**: Layer-wise activation analysis, gradient flow, architecture-specific metrics (receptive field for CNNs, memory retention for RNNs). | **Answer: Transformers**: Multi-head attention, feed-forward networks, normalization layers. **Testing**: Attention pattern analysis, layer-wise representation quality, scaling behavior, and emergent capabilities at different model sizes. |

## Question 22: How would you test different machine learning algorithms for a classification task?

| Classic AI | GenAI |
|------------|-------|
| **Answer: Algorithms**: SVM, Random Forest, Gradient Boosting, Neural Networks. **Testing**: Cross-validation, precision-recall curves, ROC analysis, feature importance, computational complexity, interpretability assessment. Compare performance across different data distributions. | **Answer: Algorithms**: Fine-tuned transformers, few-shot learning, prompt-based classification. **Testing**: Compare with traditional ML baselines, evaluate prompt sensitivity, few-shot performance, computational efficiency, and consistency across different prompt formulations. |

## Question 23: What evaluation frameworks would you use for model assessment?

| Classic AI | GenAI |
|------------|-------|
| **Answer: Frameworks**: scikit-learn metrics, MLflow, Weights & Biases, TensorBoard. **Evaluation**: Automated metrics, statistical significance testing, cross-validation, holdout testing, and benchmark datasets (UCI, Kaggle). | **Answer: Frameworks**: HuggingFace Evaluate, LangChain evaluation, OpenAI Evals, BIG-bench, HELM. **Evaluation**: Automated metrics (BLEU, ROUGE), human evaluation platforms, safety benchmarks, and domain-specific evaluations. |

## Question 24: How would you implement continuous evaluation and monitoring?

| Classic AI | GenAI |
|------------|-------|
| **Answer: Monitoring**: Model performance drift, data drift detection, prediction confidence, A/B testing infrastructure. **Tools**: MLflow, Kubeflow, custom dashboards. **Metrics**: Accuracy trends, inference latency, resource utilization. | **Answer: Monitoring**: Generation quality, safety metrics, user satisfaction, computational costs. **Tools**: LangSmith, Weights & Biases, custom evaluation pipelines. **Metrics**: Coherence scores, toxicity detection, factual accuracy, user feedback integration. |

## Question 25: Explain different types of attention mechanisms and their testing implications.

| Classic AI | GenAI |
|---|---|
| **Answer: Attention**: Additive, multiplicative, self-attention. **Testing**: Attention weight visualization, alignment quality, computational efficiency, gradient flow analysis. Validate attention focuses on relevant input regions. | **Answer: Multi-head attention**: Parallel attention computations. **Sparse attention**: Efficient long sequences. **Cross-attention**: Multi-modal alignment. **Testing**: Attention pattern analysis, head importance evaluation, scaling efficiency, and interpretability assessment. |

## Question 26: How would you test for model robustness and adversarial attacks?

| Classic AI | GenAI |
|---|---|
| **Answer: Attacks**: FGSM, PGD, C&W attacks. **Testing**: Adversarial example generation, robustness metrics, certified defenses evaluation. Test with noise injection, input perturbations, and edge cases. | **Answer: Attacks**: Prompt injection, jailbreaking, adversarial prompts. **Testing**: Red team exercises, robustness benchmarks, safety evaluations, alignment assessments. Test with manipulated inputs and edge case scenarios. |

## Question 27: What are the key considerations for testing multimodal AI systems?

| Classic AI | GenAI |
|---|---|
| **Answer: Modalities**: Vision + text, audio + text. **Testing**: Cross-modal consistency, alignment quality, missing modality handling, computational efficiency across modalities. | **Answer: Modalities**: Vision-language models, audio-text systems. **Testing**: Cross-modal understanding, generation quality, modality switching, alignment accuracy, and consistency across different input combinations. |

## Question 28: How would you evaluate the interpretability and explainability of AI models?

| Classic AI | GenAI |
|---|---|
| **Answer: Methods**: Feature importance, SHAP values, LIME, gradient-based explanations. **Testing**: Explanation consistency, human-interpretability studies, faithfulness to model behavior, and stability across similar inputs. | **Answer: Methods**: Attention visualization, probe studies, mechanistic interpretability, chain-of-thought prompting. **Testing**: Explanation quality, reasoning consistency, factual grounding, and alignment with human understanding. |

# Advanced Technical Questions

## Question 29: Explain the difference between fine-tuning and prompt engineering.

**GenAI Answer:** Fine-tuning modifies model weights through training on specific datasets, requiring computational resources and technical expertise. Prompt engineering optimizes input instructions

without changing the model, offering faster iteration but potentially less consistent results.

## Question 30: What is RLHF and why is it important?

**GenAI Answer:** Reinforcement Learning from Human Feedback trains models to align with human preferences and values. It's crucial for creating helpful, harmless, and honest AI systems by incorporating human judgment into the training process.

## Question 31: How would you test for hallucinations in language models?

**GenAI Answer:** Compare outputs against verified knowledge bases, use fact-checking APIs, implement confidence scoring, test with questions having known answers, and use human evaluation for subjective content. Monitor for consistency across similar prompts.

## Question 32: What are the different types of neural network layers and their testing considerations?

| Classic AI | GenAI |
|---|---|
| **Answer: Layers**: Dense/FC, Convolutional, Pooling, Dropout, Batch Norm, Activation. **Testing**: Layer-wise gradient analysis, activation distributions, weight initialization impact, layer ablation studies, and computational profiling. | **Answer: Layers**: Embedding, Multi-head attention, Feed-forward, Layer normalization, Positional encoding. **Testing**: Attention pattern analysis, layer-wise representation quality, normalization effectiveness, and scaling behavior. |

## Question 33: How would you test different optimization algorithms?

| Classic AI | GenAI |
|---|---|
| **Answer: Optimizers**: SGD, Adam, RMSprop, AdaGrad. **Testing**: Convergence speed, final performance, stability, hyperparameter sensitivity, computational overhead. Plot loss curves and learning rate schedules. | **Answer: Optimizers**: AdamW, Lion, Adafactor for large models. **Testing**: Training stability, memory efficiency, convergence in large-scale settings, learning rate scheduling effectiveness, and gradient clipping behavior. |

## Question 34: What are the key differences between batch, mini-batch, and online learning?

| Classic AI | GenAI |
|---|---|
| **Answer: Batch**: Full dataset per update. **Mini-batch**: Subset per update. **Online**: One sample per update. **Testing**: Convergence behavior, memory usage, computational efficiency, and final model quality across different batch sizes. | **Answer: Batch sizes**: Impact on training stability, gradient noise, memory requirements. **Testing**: Optimal batch size for different model sizes, gradient accumulation effects, and distributed training considerations. |

## Question 35: How would you test feature engineering and selection techniques?

| Classic AI | GenAI |
|---|---|
| **Answer: Techniques**: PCA, feature scaling, polynomial features, feature selection (univariate, RFE, LASSO). **Testing**: Dimensionality reduction quality, information preservation, computational efficiency, and downstream task performance. | **Answer: Feature engineering**: Tokenization strategies, vocabulary size, subword encoding. **Testing**: Tokenization quality, vocabulary coverage, out-of-vocabulary handling, and multilingual capabilities. |

## Question 36: What are ensemble methods and how would you test them?

| Classic AI | GenAI |
|---|---|
| **Answer: Methods**: Bagging, Boosting, Stacking, Voting. **Testing**: Individual model performance, diversity metrics, ensemble improvement, computational cost, and bias-variance trade-off analysis. | **Answer: Ensemble approaches**: Model averaging, mixture of experts, multi-model consensus. **Testing**: Output diversity, consistency, computational overhead, and ensemble calibration for generation tasks. |

## Question 37: How would you test regularization techniques?

| Classic AI | GenAI |
|---|---|
| **Answer: Techniques**: L1/L2 regularization, dropout, early stopping, data augmentation. **Testing**: Overfitting prevention, generalization improvement, optimal regularization strength, and computational impact. | **Answer: Techniques**: Dropout, weight decay, gradient clipping, label smoothing. **Testing**: Training stability, generalization to new domains, optimal regularization parameters, and impact on generation quality. |

## Question 38: What are the considerations for testing model compression techniques?

| Classic AI | GenAI |
|---|---|
| **Answer: Techniques**: Pruning, quantization, knowledge distillation. **Testing**: Compression ratio, performance degradation, inference speed, memory usage, and accuracy preservation across different compression levels. | **Answer: Techniques**: Model pruning, quantization, distillation for LLMs. **Testing**: Generation quality preservation, inference latency, memory footprint, and capability retention across different compression ratios. |

## Question 39: How would you test different sampling strategies during inference?

| Classic AI | GenAI |
|---|---|
| **Answer:** N/A - Sampling strategies are primarily for generative models. In classic AI, similar concept is prediction confidence thresholding or ensemble prediction aggregation. | **Answer: Strategies**: Greedy, beam search, nucleus (top-p), top-k sampling. **Testing**: Output quality, diversity, consistency, computational cost, and user preference across different sampling parameters. |

## Question 40: What are the key considerations for testing distributed training?

| Classic AI | GenAI |
|---|---|
| **Answer: Approaches**: Data parallelism, model parallelism. **Testing**: Scaling efficiency, communication overhead, synchronization issues, gradient consistency, and final model quality compared to single-node training. | **Answer: Approaches**: Data parallelism, model parallelism, pipeline parallelism. **Testing**: Training stability, gradient synchronization, memory efficiency, scaling laws, and model quality consistency across different distributed configurations. |

## Excellent Candidate (Senior Level)

- Demonstrates deep understanding of both classic AI and GenAI concepts
- Provides specific, actionable testing strategies
- Shows awareness of current challenges and solutions
- Discusses ethical considerations proactively
- Can design comprehensive test frameworks

## Good Candidate (Mid Level)

- Solid understanding of fundamental concepts
- Can explain basic testing methodologies
- Aware of common pitfalls and solutions
- Shows interest in learning new approaches

- Can work with existing test frameworks

## Needs Development (Junior Level)

- Basic knowledge of AI concepts

- Limited testing experience

- Requires guidance on methodology

- May focus on only one area (classic AI or GenAI)

- Can contribute to existing testing efforts with supervision

---

# Follow-up Questions by Experience Level

## For Senior Candidates:

- "How would you build a testing infrastructure for a GenAI system at scale?"

- "What are the unique challenges in testing multimodal AI systems?"

- "How do you balance automated testing with human evaluation?"

## For Mid-level Candidates:

- "Walk me through how you would test a chatbot for customer service"

- "What tools and frameworks have you used for AI testing?"

- "How do you prioritize which tests to run first?"

## For Junior Candidates:

- "What interests you most about AI testing?"

- "How would you approach learning about a new AI model?"

- "What do you think are the biggest challenges in AI testing?"