

Design a News Feed ML Ranking System - Meta

E7 Interview Preparation

1. Interview Overview

E7 Focus: Advanced ML concepts, technical depth, algorithm expertise

Problem Scope and Technical Complexity

- **Scale:** 2.9 billion monthly active users, 1.8 billion daily active users
- **Volume:** 4.75 billion posts per day, 100+ billion content interactions daily
- **Latency:** Sub-200ms p95 response time for feed generation
- **Throughput:** 10,000+ requests per second during peak hours

Key ML Challenges and Technical Requirements

- **Real-time ranking:** Process millions of posts in milliseconds
- **Personalization:** Individual user preferences and behavior modeling
- **Content diversity:** Balance relevance with content variety
- **Cold start:** Handle new users and new content
- **Multi-objective optimization:** Engagement, diversity, and safety

Interview Discussion Flow and Technical Deep-Dives

1. **Problem formulation** (10 minutes)
2. **Data strategy and feature engineering** (15 minutes)
3. **Model architecture and training** (20 minutes)
4. **Serving and optimization** (10 minutes)
5. **Advanced topics and edge cases** (5 minutes)

2. Problem Definition & Technical Requirements

E7 Focus: Complex ML problem formulation, advanced constraints

Problem Statement and Technical Objectives

Design a machine learning system that ranks Facebook posts in a user's news feed to maximize long-term user engagement while maintaining content diversity and safety.

Mathematical Formulation:

$$\text{Score}(u, p) = f(\theta, x_u, x_p, x_c, x_t)$$

Where:

- u : User features
- p : Post features
- c : Context features
- t : Temporal features
- θ : Model parameters

ML-Specific Success Metrics and Evaluation Criteria

- **Primary:** Time spent on platform, daily active users
- **Secondary:** Click-through rate, like rate, share rate, comment rate
- **Diversity:** Content type distribution, source diversity
- **Safety:** Harmful content detection rate

Scale Requirements and Technical Constraints

- **Users:** 2.9B monthly, 1.8B daily
- **Content:** 4.75B posts/day, 100B+ interactions/day
- **Latency:** <200ms p95, <100ms p50
- **Throughput:** 10K+ RPS peak
- **Storage:** Petabytes of user behavior data

E7 Criteria: Advanced problem decomposition, multi-objective optimization

3. Advanced Data Strategy & Feature Engineering

E7 Focus: Sophisticated feature engineering, advanced data techniques

Complex Data Sources and Collection Strategies

- **User behavior:** Clicks, likes, shares, comments, time spent
- **Content features:** Text, images, videos, metadata
- **Social graph:** Friends, pages, groups, interactions
- **Temporal data:** Time of day, day of week, recency
- **Contextual data:** Device, location, session information

Advanced Feature Engineering Techniques

User Embeddings:

- Graph neural networks for social connections
- Temporal embeddings for behavior patterns
- Multi-modal embeddings for content preferences

Content Embeddings:

- BERT-based text embeddings
- ResNet-based image embeddings
- Video understanding with 3D CNNs

Interaction Features:

- User-content interaction history
- Temporal decay functions
- Cross-feature interactions

Data Quality, Validation, and Preprocessing Pipelines

- **Real-time validation:** Feature drift detection
- **Data versioning:** Feature store with lineage tracking
- **A/B testing:** Feature flagging and gradual rollouts

- **Monitoring:** Data quality metrics and alerting

E7 Criteria: Feature store design, advanced data augmentation, data versioning

4. Advanced ML Model Design

E7 Focus: Deep understanding of model architectures, research integration

Model Selection with Detailed Technical Justification

Primary Architecture: Multi-task learning with transformer-based ranking

- **Base model:** BERT-like architecture for content understanding
- **Ranking head:** Multi-layer perceptron for score prediction
- **Auxiliary tasks:** Click prediction, engagement prediction, diversity scoring

Advanced Architecture Design and Optimization

Multi-Task Learning Framework:

$$L_{\text{total}} = \lambda_1 L_{\text{ranking}} + \lambda_2 L_{\text{click}} + \lambda_3 L_{\text{engagement}} + \lambda_4 L_{\text{diversity}}$$

Attention Mechanisms:

- Self-attention for content understanding
- Cross-attention for user-content interaction
- Temporal attention for behavior patterns

Ensemble Methods:

- Gradient boosting for feature interactions
- Deep learning for complex patterns
- Rule-based systems for safety

Multi-Task Learning, Ensemble Methods, and Advanced Techniques

- **Multi-task learning:** Joint optimization of ranking, engagement, diversity
- **Ensemble methods:** Combine multiple model types
- **Online learning:** Continuous model updates
- **Meta-learning:** Few-shot learning for new users/content

E7 Criteria: Novel architectures, research paper integration, technical innovation

5. Advanced Training & Evaluation

E7 Focus: Sophisticated training strategies, advanced evaluation methods

Advanced Training Techniques

Distributed Training:

- Parameter server architecture
- Gradient compression and quantization
- Asynchronous updates for real-time learning

Optimization Strategies:

- Adam optimizer with learning rate scheduling
- Gradient clipping for stability
- Regularization techniques (L1, L2, dropout)

Training Data Preparation:

- Negative sampling strategies
- Temporal data splitting
- Cross-validation with time series

Sophisticated Evaluation Methodology and Metrics

Offline Evaluation:

- Precision@K, Recall@K, NDCG@K
- A/B testing with statistical significance
- Counterfactual evaluation with inverse propensity scoring

Online Evaluation:

- Interleaving experiments
- Multi-armed bandit testing
- Long-term user engagement tracking

Advanced A/B Testing and Statistical Significance

- **Statistical power:** 80% power, 5% significance level
- **Sample size calculation:** Based on effect size and variance
- **Multiple testing correction:** Bonferroni correction for multiple metrics
- **Causal inference:** Instrumental variables for causal effects

E7 Criteria: Training optimization, advanced evaluation frameworks, model selection

6. Model Serving & Inference Optimization

E7 Focus: Advanced serving techniques, optimization strategies

Complex Online Inference Requirements

- **Real-time scoring:** <10ms per prediction
- **Batch processing:** Pre-compute scores for popular content
- **Caching strategies:** Redis for frequently accessed scores
- **Load balancing:** Distribute inference across multiple servers

Advanced Model Serving Architectures

Model Serving Pipeline:

1. **Feature retrieval:** Real-time feature lookup
2. **Preprocessing:** Feature normalization and transformation
3. **Inference:** Model prediction
4. **Post-processing:** Score calibration and ranking

Optimization Techniques:

- **Model quantization:** INT8 quantization for faster inference
- **Model pruning:** Remove redundant parameters
- **Knowledge distillation:** Distill large models to smaller ones
- **TensorRT optimization:** GPU acceleration

Model Compression, Quantization, and Optimization

- **Quantization:** 8-bit quantization with minimal accuracy loss
- **Pruning:** Structured and unstructured pruning
- **Distillation:** Teacher-student model training
- **Neural architecture search:** Automated model optimization

E7 Criteria: Real-time inference, model optimization, serving patterns

7. Advanced System Architecture

E7 Focus: Complex system design, advanced integration patterns

Sophisticated System Components and Interactions

Data Pipeline:

- **Stream processing:** Apache Kafka for real-time data
- **Batch processing:** Apache Spark for large-scale data
- **Feature store:** Real-time feature serving
- **Model store:** Model versioning and deployment

Serving Architecture:

- **API Gateway:** Request routing and load balancing
- **Feature Service:** Real-time feature retrieval
- **Ranking Service:** ML model inference
- **Caching Layer:** Redis for score caching

Advanced Data Flow and Processing Pipelines

Real-time Pipeline:

User Request → Feature Retrieval → Model Inference → Score Ranking → Response

Batch Pipeline:

Raw Data → Feature Engineering → Model Training → Model Deployment → A/B Testing

Complex Integration Patterns and Protocols

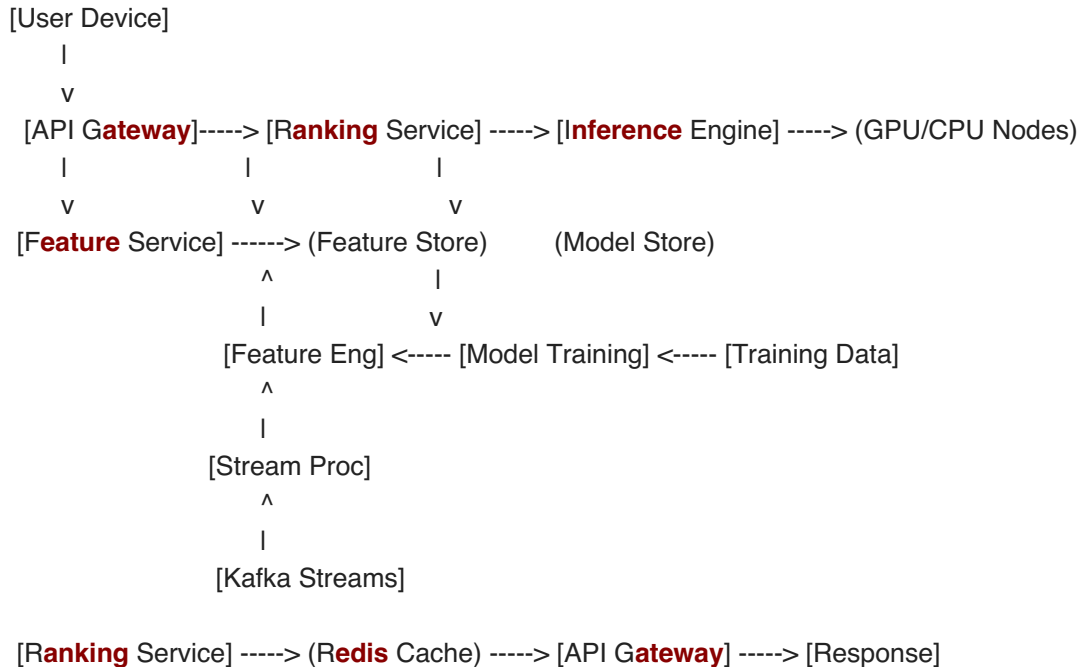
- **Microservices:** Independent scaling and deployment
- **Event-driven:** Asynchronous processing with message queues
- **Circuit breakers:** Fault tolerance and graceful degradation
- **Rate limiting:** Protect against abuse and overload

E7 Criteria: Advanced system patterns, complex data flows, optimization

System Design Diagram



ASCII Fallback Diagram



8. Advanced Monitoring & ML Observability

E7 Focus: Sophisticated ML monitoring, advanced analytics

Advanced ML Monitoring Metrics and Techniques

Model Performance Metrics:

- **Accuracy metrics:** Precision, recall, F1-score
- **Ranking metrics:** NDCG, MAP, MRR
- **Business metrics:** Engagement, retention, revenue

System Metrics:

- **Latency:** P50, P95, P99 response times
- **Throughput:** Requests per second
- **Error rates:** 4xx, 5xx error percentages
- **Resource utilization:** CPU, memory, GPU usage

Model Drift Detection and Adaptation

Drift Detection:

- **Statistical tests:** Kolmogorov-Smirnov test for distribution drift
- **Model performance:** Accuracy degradation detection
- **Feature drift:** Input distribution changes
- **Concept drift:** Label distribution changes

Adaptation Strategies:

- **Online learning:** Continuous model updates
- **Retraining:** Periodic model retraining
- **Ensemble methods:** Combine multiple models
- **Human-in-the-loop:** Expert review for edge cases

Advanced Alerting and Anomaly Detection

- **Threshold-based:** Simple rule-based alerting
- **Statistical:** Anomaly detection with statistical methods
- **ML-based:** Anomaly detection with machine learning
- **Ensemble:** Combine multiple detection methods

E7 Criteria: ML-specific monitoring, model performance analytics, drift detection

9. Advanced Scalability & Performance

E7 Focus: Complex scaling strategies, advanced optimization

Advanced Scaling Strategies for ML Workloads

Horizontal Scaling:

- **Load balancing:** Distribute requests across multiple servers
- **Sharding:** Partition data across multiple databases
- **Caching:** Multi-level caching strategy
- **CDN:** Content delivery network for global users

Vertical Scaling:

- **GPU acceleration:** CUDA for deep learning inference
- **Memory optimization:** Efficient memory usage
- **CPU optimization:** Vectorized operations
- **Storage optimization:** Efficient data storage formats

Performance Optimization and Bottleneck Analysis

Bottleneck Identification:

- **Profiling:** Identify slow components
- **Monitoring:** Real-time performance tracking
- **Load testing:** Stress testing for capacity planning
- **Optimization:** Targeted improvements

Optimization Techniques:

- **Algorithm optimization:** Efficient algorithms
- **Data structure optimization:** Efficient data structures
- **Memory optimization:** Reduce memory usage
- **Network optimization:** Reduce network latency

Resource Optimization and Efficiency

- **Auto-scaling:** Automatic resource allocation
- **Resource pooling:** Shared resources across services
- **Cost optimization:** Minimize infrastructure costs
- **Energy efficiency:** Reduce power consumption

E7 Criteria: Advanced scaling patterns, performance optimization, resource efficiency

10. Advanced ML Topics (E7 Focus)

E7 Criteria: Deep ML expertise, advanced techniques, research knowledge

Advanced ML Algorithms and Techniques

Deep Learning Architectures:

- **Transformers:** Self-attention mechanisms for sequence modeling
- **Graph Neural Networks:** Social graph analysis
- **Multi-modal Learning:** Text, image, and video understanding
- **Reinforcement Learning:** Dynamic ranking optimization

Optimization Techniques:

- **Gradient-based:** SGD, Adam, AdaGrad
- **Meta-learning:** Learning to learn
- **Neural Architecture Search:** Automated architecture design
- **Hyperparameter Optimization:** Bayesian optimization

Research Integration and Cutting-Edge Methods

Recent Research:

- **Large Language Models:** GPT, BERT, T5 for content understanding
- **Contrastive Learning:** Self-supervised learning for representations
- **Federated Learning:** Privacy-preserving distributed learning
- **Causal Inference:** Understanding cause-effect relationships

Emerging Techniques:

- **Few-shot Learning:** Learning from limited examples
- **Continual Learning:** Learning without forgetting
- **Adversarial Training:** Robust model training
- **Neural ODEs:** Continuous-time neural networks

Complex Optimization and Mathematical Foundations

Mathematical Foundations:

- **Convex Optimization:** Linear programming, quadratic programming
- **Non-convex Optimization:** Gradient descent, Newton's method
- **Stochastic Optimization:** Stochastic gradient descent
- **Multi-objective Optimization:** Pareto optimality

Advanced Techniques:

- **Bayesian Optimization:** Efficient hyperparameter tuning

- **Evolutionary Algorithms:** Genetic algorithms for optimization
- **Simulated Annealing:** Global optimization
- **Particle Swarm Optimization:** Swarm intelligence

Advanced Evaluation and Validation Techniques

Evaluation Methods:

- **Cross-validation:** K-fold, time series cross-validation
- **Bootstrap:** Resampling for confidence intervals
- **Permutation tests:** Non-parametric significance testing
- **Bayesian evaluation:** Posterior distributions

Validation Strategies:

- **Holdout validation:** Train/validation/test splits
- **Cross-validation:** Multiple train/validation splits
- **Leave-one-out:** Extreme cross-validation
- **Time series validation:** Temporal data splitting

Future ML Trends and Emerging Technologies

Emerging Trends:

- **Foundation Models:** Large pre-trained models
- **Multimodal AI:** Vision-language models
- **Causal AI:** Causal reasoning and inference
- **Neural-Symbolic AI:** Combining neural and symbolic reasoning

Technological Advances:

- **Quantum Machine Learning:** Quantum algorithms for ML
- **Edge AI:** On-device machine learning
- **Federated Learning:** Distributed privacy-preserving learning
- **AutoML:** Automated machine learning

11. Technical Deep-Dive Discussion Points

E7 Focus: Advanced technical concepts, algorithm deep-dives

Complex Technical Decisions and Trade-offs

Model Complexity vs. Performance:

- **Simple models:** Fast inference, easy to interpret
- **Complex models:** Better accuracy, slower inference
- **Trade-off:** Balance between accuracy and efficiency

Online vs. Batch Learning:

- **Online learning:** Real-time updates, adaptive to changes
- **Batch learning:** Stable training, better convergence
- **Trade-off:** Adaptability vs. stability

Accuracy vs. Diversity:

- **High accuracy:** Relevant content, user satisfaction
- **High diversity:** Content variety, user exploration
- **Trade-off:** Relevance vs. exploration

Advanced Follow-up Questions and Technical Challenges

Algorithm Questions:

- "How would you handle the cold start problem for new users?"
- "What's the computational complexity of your ranking algorithm?"
- "How would you implement online learning for your model?"

System Questions:

- "How would you handle model serving at scale?"
- "What's your strategy for A/B testing ML models?"
- "How would you detect and handle model drift?"

Research Questions:

- "How would you incorporate recent research in transformers?"
- "What's your approach to multi-modal learning?"
- "How would you implement causal inference in ranking?"

Edge Cases and Failure Scenarios

Data Issues:

- **Missing features:** Handle incomplete user data
- **Data corruption:** Detect and handle bad data
- **Feature drift:** Adapt to changing data distributions

Model Issues:

- **Model degradation:** Detect and handle performance drops
- **Overfitting:** Prevent model overfitting
- **Bias:** Detect and mitigate model bias

System Issues:

- **High latency:** Handle slow inference times
- **Service failures:** Graceful degradation
- **Resource constraints:** Handle limited resources

E7 Criteria: Algorithm complexity, mathematical foundations, advanced techniques

12. Advanced Technical Interview Questions

E7 Focus: Deep technical knowledge, advanced ML concepts

Advanced Technical Deep-Dive Questions

Algorithm Questions:

1. "Design a multi-armed bandit algorithm for content recommendation"

2. "How would you implement Thompson sampling for exploration-exploitation?"
3. "What's the difference between collaborative filtering and content-based filtering?"
4. "How would you handle the curse of dimensionality in feature engineering?"

Mathematical Questions:

1. "Derive the gradient of a neural network with respect to its parameters"
2. "Explain the mathematical foundation of attention mechanisms"
3. "How would you prove the convergence of stochastic gradient descent?"
4. "What's the relationship between regularization and bias-variance trade-off?"

System Design Questions:

1. "Design a distributed training system for large-scale neural networks"
2. "How would you implement model versioning and rollback?"
3. "Design a real-time feature store for ML systems"
4. "How would you handle model serving with zero downtime?"

Complex System Design Follow-ups

Scalability Questions:

1. "How would you scale your system to handle 10x more users?"
2. "What's your strategy for handling peak traffic?"
3. "How would you optimize for different geographic regions?"

Performance Questions:

1. "How would you reduce inference latency by 50%?"
2. "What's your strategy for memory optimization?"
3. "How would you handle model serving on mobile devices?"

Reliability Questions:

1. "How would you ensure 99.9% uptime for your ML system?"
2. "What's your disaster recovery strategy?"
3. "How would you handle data corruption in production?"

E7 Criteria: Advanced ML algorithms, mathematical depth, technical expertise

Key Takeaways for E7 Interview

1. **Technical Depth:** Demonstrate deep understanding of ML algorithms, mathematical foundations, and system design
2. **Problem Solving:** Show ability to break down complex problems and design elegant solutions
3. **Trade-offs:** Understand and articulate various trade-offs in ML system design
4. **Research Integration:** Stay current with latest research and apply it to real-world problems
5. **Scalability:** Design systems that can handle massive scale and real-time requirements
6. **Evaluation:** Implement rigorous evaluation and monitoring strategies
7. **Innovation:** Propose novel approaches and improvements to existing systems