

# Real-Time Personalization at Scale: Neural Ranking Systems and Operational Breakthroughs

Vedant Agarwal

Senior Software Engineer – Search

# Overview of the talk



- Ranking Breakthroughs
- Latency Breakthroughs
- Real-Time Constraints
- Business Impact
- Emerging Trends & Future Initiatives

# Challenges in Real-Time Personalization

---

## **Data Inconsistencies & Noise:**

- Variability in user clickstream and transaction data

## **Scaling Bottlenecks:**

- Handling millions of concurrent interactions

## **Model Drift:**

- Rapidly changing user behavior over time

## **High-Volume Streaming Data:**

- Need to process and update models continuously

# Two – Tier Neural Ranking Architecture



## **L1 – Candidate Generation:**

- Embedding-based indexing for rapid filtering

## **L2 – Precision Re-Ranking:**

- Advanced sequence modeling (LSTM/Transformer) for detailed ranking

## **Decoupled Processing:**

- Separates speed (L1) from personalization accuracy (L2)

## **Scalability:**

- Optimizes resource usage by splitting the workload

# Breakthrough in Ranking – L1: Candidate Generation

---

## **Embedding-Based Indexing:**

- Projects user behavior and item attributes into a shared space

## **Approximate Nearest Neighbor (ANN) Search:**

- Uses libraries like FAISS or Annoy for quick retrieval

## **Efficient Pre-Filtering:**

- Reduces the search space for subsequent re-ranking

## **Performance Impact:**

- High throughput with minimal latency

# Breakthrough in Ranking – L2: Precision Re-Ranking

---

## **Advanced Sequence Modeling:**

- Utilizes LSTM/Transformer models to capture sequential dependencies

## **Integration of Real-Time & Historical Data:**

- Combines session data with long-term user history

## **Attention Mechanisms:**

- Focus on the most relevant signals for each recommendation

## **Result:**

- Significantly improved recommendation accuracy

# Latency Breakthroughs & Real-Time Constraints

---

## **Sub-50ms End-to-End Processing:**

- Ensures rapid response from user action to recommendation

## **Model Complexity vs. Inference Speed:**

- Simplified architecture to meet time constraints without sacrificing quality

## **Optimized Inference Pipelines:**

- Asynchronous processing and caching strategies in place

## **Scalability Under Load:**

- Dynamic allocation (autoscaling)

# Infrastructure & Feature Engineering

---

## **Real-Time Data Pipelines:**

- Utilizes Kafka and Apache Flink for continuous data ingestion

## **Microservices Architecture:**

- Containerized services via Docker and orchestrated by Kubernetes

## **Dynamic Feature Engineering:**

- Generates features on-the-fly from live data (e.g., user session attributes)
- Advanced transformations: normalization, feature crosses, deep embeddings

## **Observability:**

- Monitoring using Prometheus and Grafana for real-time debugging



# Deployment Strategies & Optimization Practices

---

## **Regular Monitoring & KPIs:**

- Automated anomaly alerts to catch performance deviations early

## **Automated Retraining:**

- Scheduled model updates to address drift and ensure freshness

## **Blue-Green Deployments & Rolling Updates:**

- Ensures zero downtime during feature rollouts

## **A/B Testing:**

- Gradual feature rollouts with real-time impact assessment

# Business Impact & Measurable Outcomes



## **Increased Conversion Rates:**

- Direct improvement in matching products to user intent

## **Enhanced User Engagement:**

- Personalized recommendations drive higher CTR and longer sessions

## **Improved Customer Retention:**

- Context-aware suggestions lead to repeat interactions

# Emerging Trends & Future Initiatives

---

## **Multimodal Personalization:**

- Integrates text, image, and video data for richer insights

## **Real-Time Federated Learning:**

- Decentralized model training to enhance privacy and reduce latency

## **Enhancing Neural Ranking:**

- Adoption of next-gen deep learning models for further accuracy gains

## **Hybrid Approaches:**

- Combining rule-based systems with AI for improved interpretability

## **Long-Term Vision:**

- Continuously adapt to market trends and evolving user behaviors

# Conclusion

---

## **Recap of Breakthroughs:**

- Ranking innovations (L1 & L2), latency optimizations, real-time scalability

## **Business Outcomes:**

- Measurable improvements in conversion, engagement, and retention

## **Future Outlook:**

- Emerging trends and next-generation personalization initiatives

Thank You!