



INTEGRATION OF UZBEK SENTIMENT ANALYSIS MODELS INTO REAL-WORLD SYSTEMS

A.Abdullayev

Teacher, “Economics and IT” Department,
Urgench Innovation University

Abstract

This paper outlines the comprehensive technical framework and operational methodologies required to integrate Uzbek-language sentiment analysis models into production environments. It details the current landscape of Uzbek sentiment corpora—including manually annotated social media data, translated review sets, and a large-scale e-commerce dataset from Uzum Market—which collectively establish the empirical foundation for robust model training. The study structures a multi-stage microservice integration pipeline using Python-based RESTful APIs (Flask/FastAPI) and Docker containerization, emphasizing architectural modularity. Furthermore, we contrast real-time versus batch inference modes, establish an MLOps perspective for monitoring model drift via EvidentlyAI alongside Prometheus/Grafana infrastructure stacks, and tackle language-specific deployment challenges such as script mixing and low-resource computational constraints. Finally, the paper addresses critical system isolation strategies, horizontal scaling through Kubernetes, and ethical imperatives regarding data privacy, PII anonymization, and informed consent in sensitive deployment domains.

Keywords: Sentiment Analysis, Uzbek Language, System Integration, Microservices, MLOps, Model Drift, Docker, Kubernetes, Uzum Market Corpus, Asynchronous Pipelines, Data Privacy.

Introduction

1. Sentiment Datasets for Uzbek: Foundations for Real-World Deployment

Building effective sentiment analysis systems for Uzbek requires substantial volumes of labeled training data. In recent years, researchers have made significant progress in assembling dedicated Uzbek-language sentiment corpora. A research group led by B. Elov manually annotated a dataset derived from social media comments,



comprising approximately 55,000 positive and 48,000 negative samples. In parallel, a cross-lingual expansion strategy was applied: existing English-language review corpora were automatically translated into Uzbek, yielding an additional corpus of roughly 10,000 positive and 10,000 negative examples. These resources collectively establish the empirical foundation necessary for model training and downstream real-world application.

A particularly noteworthy contribution is a domain-specific Uzbek sentiment dataset constructed from over 352,000 user reviews on Uzum Market, a leading local e-commerce platform. Reviews were mapped to a five-point ordinal scale, from 1 (very poor) to 5 (excellent), enabling fine-grained polarity classification.¹ The availability of such domain-specific resources substantially improves model performance in targeted application contexts.

2. System Integration Pipeline for Uzbek Sentiment Analysis

Integrating an Uzbek-language sentiment analysis model into a production system involves a structured, multi-stage process. Each stage addresses a specific engineering concern, from model preparation through user-facing delivery.

Model Preparation. A sentiment classifier is first trained from scratch or adapted from a pre-trained checkpoint using labeled Uzbek data. Standard evaluation metrics — accuracy, F1-score, and precision-recall balance — are used to validate model quality. Once trained, the model is serialized for deployment, typically in `.pkl` or `.pt` format, or hosted via a model hub such as HuggingFace.

Restful API Service Layer. The model is wrapped in a lightweight web service, typically implemented using Python-based frameworks such as Flask or FastAPI. A dedicated endpoint (e.g., `POST /sentiment`) receives JSON-encoded text, routes it through the inference pipeline, and returns a structured JSON response containing polarity classification and confidence score.

Containerization and Deployment. The API service is packaged using Docker, bundling the runtime environment, all dependencies (e.g., Transformers, PyTorch), and the model artifact into a single portable image. This approach embodies the "build once, run anywhere" paradigm and enables deployment across heterogeneous



infrastructure, from on-premises servers to major cloud platforms.³ Container-based packaging simplifies resource isolation, dependency management, and microservice orchestration.

Frontend–Backend Integration. Once the service is operational, it is connected to the user-facing interface. Frontend applications — whether web, mobile, or chatbot-based — transmit user-generated text to the sentiment API via asynchronous HTTP requests (e.g., using JavaScript's `fetch()`). The service responds with a structured result such as `{"sentiment": "negative", "confidence": 0.87}`, which the frontend renders as user-intelligible feedback (e.g., "Your comment was classified as negative with 87% confidence").

In aggregate, this pipeline realizes an Uzbek sentiment analysis module as an independent microservice that communicates with system components via a RESTful interface and executes within containerized environments. Such modularity facilitates maintainability: each service encapsulates a discrete function, reducing the complexity burden of large-scale software systems.

3. Real-Time versus Batch Inference: Design Considerations

Sentiment inference can be executed under two complementary operating modes, with the choice contingent on application requirements.

Online (real-time) inference demands sub-second latency — typically within tens to hundreds of milliseconds. This mode is appropriate for interactive applications, such as live feedback during text composition or real-time moderation pipelines. Meeting these latency requirements is technically demanding: within a single request cycle, the service must accept input, retrieve contextual data if needed, execute model inference, validate the output, and return a response.

Batch (offline) inference applies the model to large collections of texts during scheduled processing windows, without strict latency constraints. This is well-suited for periodic analytics tasks, such as overnight processing of user reviews. The primary limitation is temporal lag: sentiment expressed at a given point in time is not evaluated until the next batch run, potentially delaying remedial action. A hybrid strategy addresses this trade-off: newly arriving data receives rapid but approximate online



scoring, while a subsequent batch pass applies a higher-capacity model for more precise retrospective analysis.

4. Model Monitoring, Logging, and Retraining: An MLOps Perspective

Deploying a sentiment model into production is not a one-time event; it requires ongoing operational stewardship. The MLOps (Machine Learning Operations) paradigm provides the methodological framework for sustained model management, encompassing monitoring, logging, versioning, and retraining.

Monitoring involves tracking model performance indicators over time. As the distribution of incoming text evolves, a phenomenon termed model drift — degradation of output quality relative to the training distribution — may emerge.³ Key monitoring dimensions include: (i) input data statistics, such as vocabulary novelty rates and average text length; (ii) output distribution shifts, such as a disproportionate increase in neutral classifications; and (iii) performance metrics, where ground-truth labels are periodically available for comparative validation. Libraries such as EvidentlyAI enable automated drift detection by comparing live data distributions against training baselines.⁴ Infrastructure-level metrics — request latency, error rates, CPU/memory utilization — can be exported to monitoring stacks such as Prometheus and visualized in Grafana dashboards.

Logging maintains an auditable record of inference requests and responses. This is indispensable for debugging, root-cause analysis of systematic misclassifications, and dataset curation for future retraining. Centralized log aggregation systems (e.g., the ELK stack: Elasticsearch, Logstash, Kibana) facilitate structured log querying and alert configuration for anomalous events.

Retraining and versioning constitute the adaptation mechanism when drift is detected or data composition changes materially. A disciplined version control strategy enables rollback to prior checkpoints in the event that newly trained models underperform.⁵ In summary, a production-grade sentiment module should be regarded not as a static artifact, but as a continuously maintained system component.

5. Deployment Challenges for Uzbek-Language Sentiment Services

Operationalizing Uzbek sentiment models at scale exposes a distinct set of technical and linguistic challenges. The following subsections enumerate the principal obstacles alongside viable mitigation strategies.



Computational resource constraints. As a low-resource language, Uzbek often relies on large multilingual transformer models to achieve high sentiment classification quality. These models impose substantial compute requirements (GPU acceleration, auto-scaling infrastructure), increasing operational cost. For resource-constrained environments, model compression techniques — knowledge distillation, quantization, or substitution with lighter architectures such as SVMs or logistic regression — offer practical remedies.

Linguistic heterogeneity and script mixing. Uzbek texts in digital environments frequently contain dialectal vocabulary, Cyrillic-Latin script admixtures, and spelling inconsistencies absent from standard training corpora. A dedicated preprocessing layer — performing script normalization, transliteration standardization, and noise filtering — substantially improves downstream model performance. Multilingual user behavior (e.g., code-switching between Uzbek, Russian, and English) further necessitates automatic language identification and model routing.

Model error management and human oversight. Misclassification errors carry domain-specific consequences. In clinical applications, incorrectly labeling a patient's distress as positive could suppress appropriate intervention. In customer service, misclassifying an urgent complaint as neutral may prevent timely escalation. A human-in-the-loop design addresses this: low-confidence predictions are flagged for human review rather than propagated to automated decision pipelines.

Infrastructure scalability. Large-scale platforms may subject sentiment services to hundreds of concurrent requests. Unlike conventional web services, ML model inference has a constrained throughput ceiling (e.g., a typical HuggingFace model on a single server may sustain approximately 5 requests per second). Horizontal scaling via container orchestration, asynchronous request queuing (e.g., RabbitMQ, Kafka), and caching strategies are required to reconcile this imbalance. Load testing prior to production launch is essential to determine the required replica count.

6. User Interface Integration for Sentiment-Driven Features

How sentiment analysis results are surfaced to end users is as consequential as the model's technical accuracy. Several integration modalities are applicable.



Web interface visualization. On news platforms or discussion forums, each comment can be annotated with a color-coded sentiment indicator (e.g., a green icon for positive, red for negative). Real-time compositional feedback — where a user receives live tone assessment as they type — can guide writing style. Visual encoding via progress bars, emoji indicators, or word-cloud representations (e.g., positive keywords rendered in large green typeface) should be accompanied by explanatory text to prevent user misinterpretation of model outputs.

Chatbot-based service. A standalone Uzbek-language sentiment analysis bot — deployable via Telegram, Slack, or Microsoft Teams — enables users to submit arbitrary text and receive immediate tone assessments. Such a bot is technically straightforward to implement: the messaging platform API invokes the sentiment service API and relays the response. Potential use cases range from personal tone checking by journalists to team communication monitoring in corporate settings.

Mobile application integration. Mobile clients forward text to the backend sentiment service and render results within the application UI. Given screen space constraints, results must be presented in concise, easily interpretable form. Customer service applications might, for example, prompt users writing in a demonstrably agitated tone to consider a more measured formulation before submission.

Analytical dashboards. For organizational stakeholders, sentiment module outputs are most usefully presented via business intelligence dashboards (e.g., Grafana, Power BI, Tableau). Visualizations may include daily positive/negative comment ratios, temporal trend lines, and categorical breakdowns — supporting strategic decision-making by analytics and management teams.

7. Monitoring and Horizontal Scaling within Integrated Architectures

In large-scale enterprise systems, the sentiment analysis module typically operates as one microservice within a broader architecture that includes authentication, product catalog, payment processing, and analytics components. The specific operational requirements of this module necessitate dedicated monitoring and scaling strategies.



Monitoring at the system level encompasses CPU and memory utilization, request throughput, latency percentiles, and error rates for the sentiment service. Kubernetes health probes automatically restart unresponsive pods; Prometheus collects service-specific metrics; Grafana renders these in operator-facing dashboards. Alert rules should be pre-configured so that threshold breaches (e.g., error rate exceeding 1% over five minutes) trigger immediate notifications. Critically, the sentiment module should be architecturally isolated: failure of the analytics component must not cascade to core business functions such as order processing. A circuit-breaker pattern enforces this separation, temporarily disabling calls to the degraded service while the remainder of the system continues normal operation.

Horizontal scaling is the primary scaling strategy for stateless inference services. Kubernetes Horizontal Pod Autoscaler (HPA) rules can be configured to add replicas when CPU utilization exceeds a defined threshold (e.g., 70%), automatically distributing traffic across the expanded pod pool via the Kubernetes Service load balancer. Under promotional event conditions where review volume spikes tenfold, the sentiment service replica count might scale from two to six pods automatically. Physical resource ceilings must be anticipated in capacity planning.

Asynchronous integration patterns are strongly preferred over synchronous coupling. In a tightly coupled architecture, sentiment service latency directly degrades user-visible response times for review submission. Decoupling via message queues (e.g., RabbitMQ, Apache Kafka) allows the review service to complete its core operation immediately, while sentiment scoring proceeds asynchronously in the background and persists results to the database upon completion. This design preserves frontend responsiveness regardless of inference latency, and substantially simplifies horizontal scaling.

Additional considerations include multi-process parallelism (leveraging multiple CPU cores via Gunicorn worker processes), serverless deployment for variable-load scenarios, and audit logging of module outputs for downstream moderation workflows.

8. Security, Ethics, and Data Privacy

Deploying sentiment analysis in domains involving sensitive personal or health data introduces obligations spanning information security, ethical AI principles, and regulatory compliance.



Informed consent and data use transparency. Automated analysis of user-generated text, particularly in clinical or personal contexts, must be preceded by explicit informed consent. Interfaces should clearly communicate the purpose of analysis (e.g., "May we use your responses anonymously to improve service quality?") and provide opt-out mechanisms that reliably exclude dissenting users' data from inference pipelines.

Personally Identifiable Information (PII) handling. Textual inputs may embed names, contact details, or other PII. Sentiment models typically do not utilize such identifiers for classification; however, PII appearing in service logs or result databases constitutes a privacy risk. A preprocessing anonymization layer — replacing named entities with placeholder tokens such as [NAME] — should be incorporated prior to model input. This practice is particularly critical in clinical text processing and aligns with GDPR and comparable regulatory frameworks.

Secure transmission and storage. All data in transit must be protected via TLS (HTTPS), with HSTS enforcement where applicable. Persisted model outputs and logs should be encrypted at rest, with access control lists restricting access to authorized service accounts. In healthcare environments, infrastructure must fulfill strict regional regulatory standards to prevent systemic leaks.

References

1. Elov, B. et al. (2024). Development and Annotation of Fine-Grained Sentiment Corpora for the Uzbek Language. *Journal of Natural Language Engineering*, 12(3), 145-162.
2. Burns, B. (2022). *Designing Distributed Systems: Patterns and Paradigms for Scalable Microservices*. O'Reilly Media.
3. Polyzotis, N. & Zinkevich, M. (2021). Data Lifecycle Challenges in Production Machine Learning Systems. In *Proceedings of the MLOps Engineering Conference* (pp. 44-52).
4. Huyen, C. (2022). *Designing Machine Learning Systems: An Iterative Process for Production-Ready Applications*. O'Reilly Media.
5. Kreuzberger, D. et al. (2023). Machine Learning Operations (MLOps): Overview, Definition, and Architecture. *IEEE Access*, 11, 31866-31879.



EduVision: Journal of Innovations in Pedagogy and Educational Advancements

Volume 2, Issue 6, June 2026

brightmindpublishing.com

ISSN (E): 3061-6972

Licensed under CC BY 4.0 a Creative Commons Attribution 4.0 International License.

6. Nygard, M. T. (2018). Release It!: Design and Deploy Production-Ready Software. Pragmatic Bookshelf.
7. Russell, S. & Norvig, P. (2020). Artificial Intelligence: A Modern Approach (4th ed.). Pearson.