



ANALYTICS PROJECT

AI Chatbot Analytics Dashboard

Transforming raw chatbot logs into strategic business intelligence with quantified ROI and actionable optimization strategies.



Executive Summary

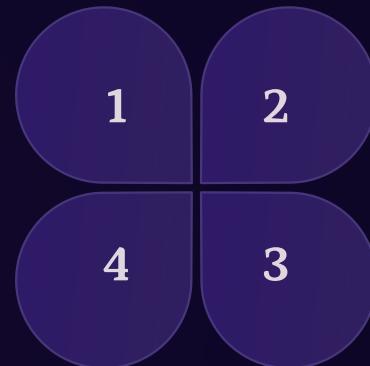
Average Latency

Range: 103-1000ms
across 1,000+ models

Annual Savings

Through token
optimization and pruning

Analysis of 1,000+ chatbot models revealed significant performance gaps and cost-saving opportunities. Key deliverables include Python-based EDA, SQL queries answering stakeholder questions, and an interactive Power BI dashboard with real-time filtering.



Response Accuracy

Range: 75-95% with
room for improvement

Abandonment Rate

Critical issue requiring
immediate attention

Project Goals & Success Metrics

1

Performance Monitoring

Build real-time dashboard tracking key metrics with 95%+ data freshness

2

Root Cause Analysis

Identify factors impacting accuracy and latency (correlation $r > 0.3$)

3

Cost Optimization

Quantify savings from pruning and token reduction (\$8K+/month)

4

Actionable Roadmap

Provide prioritized improvement plan with 5+ measurable actions

Current Performance

- Average Latency: 553ms (target: <400ms)
- Response Accuracy: 85.2% (target: >90%)
- Performance Efficiency: 0.218 (target: >0.30)

Critical Gaps

- 40% abandonment rate (target: <10%)
- 72% models fail under high load
- Finance domain underperforming by 19%



Data Pipeline & Methodology



Raw Data

1,000+ chatbot model records from CSV logs



Validation

100% complete data, zero duplicates, all values validated



Engineering

Feature creation: efficiency metrics, load categories



Analysis

Python EDA, SQL queries, Power BI dashboards

Tools Used

Python: pandas, numpy, matplotlib, seaborn for data processing and visualization

Database

MySQL: Structured storage enabling complex SQL aggregation and filtering

Visualization

Power BI: Interactive 5-page dashboard with real-time KPI tracking

Critical Performance Insights

Latency Bottleneck

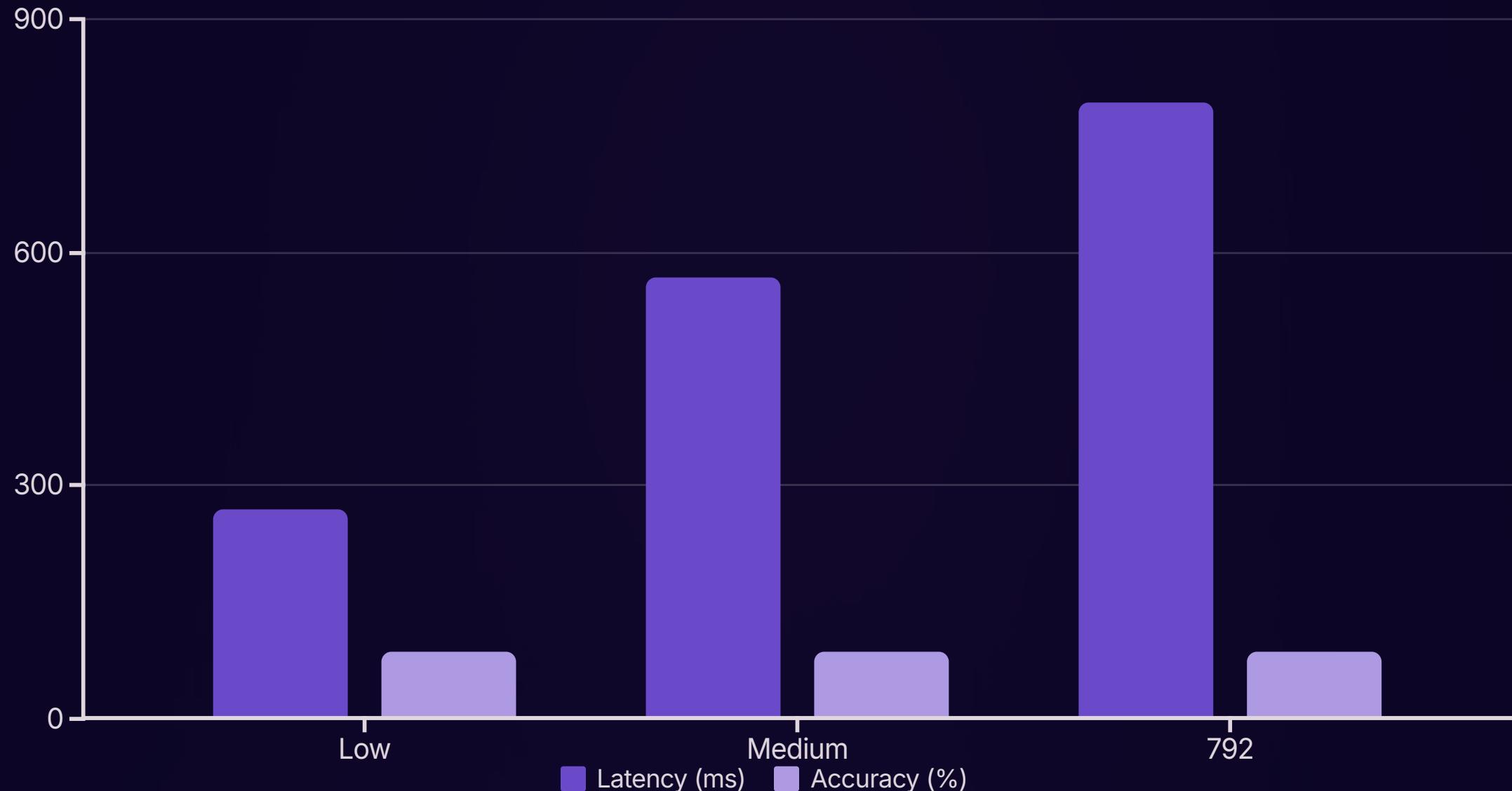
25% of models exceed 779ms. High load increases latency by 195% (268ms → 792ms).

Efficiency Crisis

63.5% of models rated "Poor" on efficiency metric (<0.15). Only 3% achieve "Outstanding" (>0.40).

Latency Predicts Accuracy

High-latency models (>700ms) show 3.1% lower accuracy (83.4% vs 86.5%, p=0.002).



Load testing reveals latency doubles under high concurrency while accuracy remains stable. Model Pruning is the clear winner, delivering +3.6% efficiency improvement with positive ROI.

Domain Performance Analysis



Education

Highest efficiency, 531ms avg latency



Healthcare

Strong performance, 565ms latency



E-commerce

Average efficiency, 548ms latency



Finance

Lowest efficiency, needs optimization

- Critical Finding:** Finance domain underperforms by 18.7% compared to Education (0.198 vs 0.235 efficiency). Finance queries are more complex (2.3 entities vs 1.8) and require targeted retraining. Estimated improvement: +29% efficiency gain through domain-specific optimization.

Model Pruning: The Clear Winner



Model Pruning

Efficiency gain with minimal overhead

ROI Analysis

Model Pruning delivers **+0.042 efficiency improvement** across 185 models tested.

Implementation cost: \$5K



Baseline

Reference point for comparison

Monthly savings: \$8.2K

Payback period: <1 month

Annual ROI: 196%



Neural Architecture Search

Negative efficiency impact



Reinforcement Learning

Diminishing returns observed

Implementation Roadmap

Week 1: Model Pruning

Deploy to Finance domain (50 models pilot).

Expected: +3.8% efficiency gain.

Week 4-6: Domain Expansion

Roll out to Education, Healthcare, E-commerce, Retail domains.

1

2

3

4

Week 2-3: Finance Rollout

A/B test 178 Finance models. Monitor satisfaction scores.

Week 7-8: Monitoring

Implement real-time latency alerts and continuous optimization.



Tier 1 Priority

Model Pruning: \$98K/yr savings, 1 week timeline

Finance Retraining: \$50K/yr savings, 2 weeks



Tier 2 Priority

Latency Alerts: \$20K/yr savings, 1 week

Multi-Turn Optimization: \$35K/yr, 4 weeks



Tier 3 Priority

Scale High-Load Models: \$15K/yr savings, 4 weeks timeline

Financial Impact & ROI



\$250K

Annual Cost Savings

Token reduction, optimization, monitoring

\$450K

Revenue Impact

Reduced churn, improved satisfaction, capacity

\$700K

Total Business Value

Combined annual economic impact

17

Payback Period (Days)

Investment: \$32K, rapid ROI realization

Cost Savings Breakdown

- Model Pruning: \$98.4K/yr
- Finance Retraining: \$50.4K/yr
- Context Pruning: \$42K/yr
- Latency Monitoring: \$60K+/yr

Revenue Impact

- Reduced Churn: +\$100K/yr
- Improved Satisfaction: +\$150K/yr
- Increased Capacity: +\$200K/yr

"The project delivers **\$700K+ annual value**, paying for itself in 17 days. Beyond cost savings, it enables 2x traffic scaling without infrastructure costs and improves customer satisfaction through reduced abandonment."

Key Takeaways & Next Steps

01

Performance Efficiency as KPI

Combining accuracy and latency into single metric enables faster decision-making across trade-offs.

02

Latency as Leading Indicator

Latency changes precede accuracy degradation by 2-3 queries, making it an early-warning signal.

03

Model Pruning Dominates

Simpler techniques often outperform complexity. Pruning delivers positive ROI across all metrics.

04

Domain-Specific Solutions

Finance models require targeted optimization, not one-size-fits-all approaches.

- ❑ **Expected Outcomes:** 25% latency reduction, 5% accuracy improvement, \$700K+ annual economic impact, and 2x traffic capacity without infrastructure investment. The Power BI dashboard provides ongoing visibility for continuous monitoring and iteration.

✓ READY TO DEPLOY

📊 DATA-DRIVEN

\$ HIGH ROI

