

---

 **AGENDA DA APRESENTAÇÃO - SentiBR**

---

---

 **ESTRUTURA COMPLETA DA APRESENTAÇÃO**

---

---

 **ABERTURA & CONTEXTO (3 min)**

- |— 1.1 Visão Geral do Desafio
  - |— 1.2 Proposta de Solução: SentiBR Platform
  - |— 1.3 Arquitetura High-Level
  - |— 1.4 Tech Stack Overview
- 

---

 **DATA PIPELINE & FEATURE ENGINEERING (4 min)**

- |— 2.1 Aquisição e Exploração de Dados
  - |— Dataset: Brazilian Portuguese Restaurant Reviews
  - |— EDA (Exploratory Data Analysis)
  - |— Data Quality Assessment
  - |— Statistical Profiling
- |— 2.2 Data Preprocessing Pipeline
  - |— Text Normalization & Cleaning
  - |— Tokenization Strategy
  - |— Train/Val/Test Split (70/15/15)
  - |— Data Augmentation Techniques
- |— 2.3 Feature Store & Versionamento
  - |— DVC (Data Version Control) - Ready
  - |— Feature Engineering
  - |— Data Validation (Great Expectations)

## └ 2.4 Synthetic Data Generation

  └ Fallback Strategy para Dataset Issues

## └ Quality Assurance

---

## § MODEL TRAINING & FINE-TUNING (5 min)

### └ 3.1 Model Selection & Architecture

  |  └ Base Model: neuralmind/bert-base-portuguese-cased

  |  └ Transfer Learning Strategy

  |  └ Architecture Customization

  |  └ Model Complexity vs Performance Trade-off

  |

### └ 3.2 Training Pipeline

  |  └ Hyperparameter Configuration

  |  |  └ Learning Rate: 2e-5

  |  |  └ Batch Size: 16

  |  |  └ Epochs: 3-4

  |  |  └ Max Sequence Length: 128

  |  |

  |  └ Optimization Strategy

  |  |  └ AdamW Optimizer

  |  |  └ Linear LR Scheduler with Warmup

  |  |  └ Gradient Accumulation

  |  |

  |  └ Regularization Techniques

  |  |  └ Dropout (0.1)

  |  |  └ Weight Decay

  |  |  └ Early Stopping

  |  |

- |   └ Training Monitoring
  - |     ├ Loss Curves (Train/Val)
  - |     ├ Learning Rate Schedule
  - |     └ Gradient Norms
- |
- |   └ 3.3 Hyperparameter Tuning
  - |     ├ Optuna Integration (Ready)
  - |     ├ Search Space Definition
  - |     ├ Multi-objective Optimization
  - |     └ Best Configuration Selection
- |
- └ 3.4 Experiment Tracking & MLflow
  - ├ Experiment Logging
  - ├ Artifact Management
  - ├ Model Registry
  - └ Reproducibility Guarantees

---

## 4 MODEL VALIDATION & EVAL FRAMEWORK (4 min)

- └ 4.1 Evaluation Metrics Suite
  - |   └ Classification Metrics
    - |     └ Accuracy: 92.3%
    - |     └ Precision: 0.93
    - |     └ Recall: 0.90
    - |     └ F1-Score: 0.91
    - |     └ ROC-AUC: 0.96
  - |
  - |   └ Confusion Matrix Analysis
  - |   └ Per-Class Performance Breakdown

| └ Calibration Metrics (Brier Score)

|

| └ 4.2 Cross-Validation Strategy

| | └ Stratified K-Fold

| | └ Bootstrap Confidence Intervals

| | └ Statistical Significance Tests

|

| └ 4.3 Error Analysis & Debugging

| | └ Failure Case Investigation

| | └ Bias Detection

| | └ Edge Cases Identification

| | └ Troubleshooting Guidelines

|

└ 4.4 LLM-as-Judge EVAL Framework ★

| └ GPT-4o-mini Validation Layer

| └ Agreement Score (Cohen's Kappa)

| └ Disagreement Analysis

└ Quality Assurance Pipeline

---

## 5 MODEL IMPROVEMENT & OPTIMIZATION (3 min)

| └ 5.1 Performance Optimization

| | └ Model Quantization (INT8)

| | └ ONNX Conversion (Ready)

| | └ TorchScript Compilation

| | └ Latency Profiling

|

| └ 5.2 Continuous Learning Loop

| | └ Feedback Collection System

- |   └ Active Learning Strategy
- |   └ Incremental Training Pipeline
- |   └ Model Versioning & Rollback
- |
- |   └ 5.3 Ensemble & Hybrid Approaches
- |   |   └ BERT + GPT Ensemble (Planned)
- |   |   └ Weighted Averaging
- |   |   └ Confidence-based Routing
- |   |
- |   └ 5.4 A/B Testing Framework
- |   └ Traffic Splitting Strategy
- |   └ Statistical Testing
- |   └ Winner Selection Criteria

---

## **6 BACKEND API & INFERENCE ENGINE (4 min)**

- |   └ 6.1 API Architecture (FastAPI)
  - |   |   └ RESTful Endpoints (11 total)
    - |   |   |   └ POST /predict - Single Inference
    - |   |   |   └ POST /predict/compare - BERT vs GPT ★
    - |   |   |   └ POST /predict/batch - Batch Processing
    - |   |   |   └ POST /feedback - Feedback Loop
    - |   |   |   └ GET /metrics - Prometheus Export
    - |   |   |   └ GET /health - Health Checks
    - |   |   |   └ GET /explain - Explainability
    - |   |   |
    - |   |   └ Async/Await Pattern
    - |   |   └ Request Validation (Pydantic)
    - |   |   └ Error Handling & HTTP Status Codes

- |
  - |— 6.2 Inference Optimization
    - |— Model Loading Strategy
    - |— Batch Inference (Dynamic Batching)
    - |— GPU Acceleration (CUDA Support)
    - |— Multi-threading/Multi-processing
  - |
  - |— 6.3 Caching Layer (Redis)
    - |— Response Caching
    - |— TTL Strategy
    - |— Cache Invalidation
    - |— Hit Rate Monitoring
  - |
  - |— 6.4 Database Layer (PostgreSQL)
    - |— Prediction Logging
    - |— Feedback Storage
    - |— Audit Trail
    - |— Query Optimization
  - |
  - |— 6.5 Rate Limiting & Security
    - |— API Key Authentication (Ready)
    - |— Rate Limiting (per-user)
    - |— CORS Configuration
    - |— Input Sanitization
  - |
  - |— 6.6 Explainability Engine
    - |— LIME Integration
    - |— SHAP Values (Ready)

- |— Attention Visualization
  - |— Feature Importance Ranking
- 

## **OBSERVABILITY & TELEMETRY (4 min)**

- |— 7.1 Logging Infrastructure
  - |— Structured Logging (JSON Format)
  - |— Log Levels (DEBUG/INFO/WARNING/ERROR)
  - |— Contextual Information
  - |— Log Aggregation Strategy
- |— 7.2 Metrics Collection (Prometheus)
  - |— System Metrics
    - |— CPU/Memory Usage
    - |— Disk I/O
    - |— Network Traffic
  - |— Application Metrics
    - |— Request Rate (req/min)
    - |— Latency Percentiles (P50/P95/P99)
    - |— Error Rate
    - |— Throughput
  - |— ML Model Metrics
    - |— Prediction Distribution
    - |— Confidence Scores
    - |— Model Load Time
    - |— Inference Time

- |   └ Business Metrics
    - |     ├ Predictions per Sentiment Class
    - |     ├ User Feedback Rate
    - |     └ API Usage Patterns
  - |
  - |   └ 7.3 Distributed Tracing (Ready)
    - |     ├ Request ID Propagation
    - |     ├ Span Collection
    - |     ├ Trace Visualization
    - |     └ Performance Bottleneck Detection
  - |
  - |   └ 7.4 Health Checks & Probes
    - |     └ Liveness Probe (API Up/Down)
    - |     └ Readiness Probe (Model Loaded)
    - |     └ Dependency Health (DB, Redis, etc)
    - |     └ Startup Probe
- 

## 8 MONITORING & ALERTING (4 min)

- |   └ 8.1 Grafana Dashboards (3 dashboards)
  - |     └ System Health Dashboard
    - |       └ Service Uptime
    - |       └ Resource Utilization
    - |       └ Request Volume
  - |     |
  - |     └ Model Performance Dashboard
    - |       └ Prediction Trends
    - |       └ Confidence Distribution
    - |       └ Error Rate by Endpoint

- | | └ BERT vs GPT Comparison
- | |
- | └ Business Intelligence Dashboard
  - | ├ Usage Heatmaps
  - | ├ Sentiment Distribution
  - | └ Feedback Analytics
- |
- | └ 8.2 Alerting System (15+ alerts)
  - | ├ Performance Alerts
    - | | └ High Latency (P95 > 200ms)
    - | | └ Low Throughput
    - | | └ Resource Exhaustion
  - | |
  - | | └ Error Alerts
    - | | | └ High Error Rate (> 5%)
    - | | | └ Model Load Failures
    - | | | └ Database Connection Issues
  - | | |
  - | | └ ML-Specific Alerts
    - | | | └ Data Drift Detection ★
    - | | | └ Prediction Drift
      - | | | | └ Low Confidence Rate (> 30%)
      - | | | | └ Anomalous Prediction Patterns
    - | | | |
    - | | | └ Business Alerts
      - | | | | └ Unusual Traffic Patterns
      - | | | | └ Negative Feedback Spike
  - | | |

## └ 8.3 Data Drift Detection Pipeline ★

- └ Statistical Tests
  - |  └ Kolmogorov-Smirnov Test
  - |  └ Chi-Square Test
  - |  └ Population Stability Index (PSI)
  - |
- └ Feature Distribution Monitoring
- └ Baseline Comparison
- └ Automated Alerts (Threshold: 0.05)
- └ Drift Visualization

---

## 9 TESTING INFRASTRUCTURE (3 min)

- └ 9.1 Unit Testing Suite
  - |  └ Model Components Tests (pytest)
  - |  └ API Endpoint Tests
  - |  └ Data Pipeline Tests
  - |  └ Utility Functions Tests
  - |  └ Coverage: 78%
  - |
- └ 9.2 Integration Testing
  - |  └ End-to-End Workflows
  - |  └ API → Model → Response
  - |  └ Feedback Loop Testing
  - |  └ Database Integration Tests
  - |  └ External Dependencies Mocking
  - |
- └ 9.3 Load Testing (Locust)
  - |  └ Performance Benchmarking

- | |   └ 100 concurrent users
- | |   └ 1000+ req/min throughput
- | |   └ P95 latency < 100ms
- | |
- |   └ Stress Testing
- |   └ Spike Testing
- |   └ Endurance Testing
- |
- └ 9.4 Model Testing & Validation
  - |   └ Regression Tests (performance baseline)
  - |   └ Adversarial Testing
  - |   └ Bias Testing
  - |   └ Edge Case Testing
- |
- └ 9.5 CI/CD Pipeline (GitHub Actions)
  - └ Automated Testing on PR
  - └ Code Quality Checks (Black, Flake8)
  - └ Docker Image Build
  - └ Deployment Automation (Ready)

---

## 10 CONTAINERIZATION & ORCHESTRATION (3 min)

- └ 10.1 Docker Architecture (8 services)
  - |   └ sentibr-api (FastAPI Backend)
  - |   └ sentibr-frontend (Streamlit)
  - |   └ postgres (Database)
  - |   └ redis (Cache)
  - |   └ prometheus (Metrics)
  - |   └ grafana (Visualization)

- |   └ mlflow (Experiment Tracking)
- |   └ nginx (Load Balancer - Ready)
- |
- |
- | └ 10.2 Docker Compose Orchestration
- |   └ Service Dependencies
- |   └ Network Configuration
- |   └ Volume Management (Persistence)
- |   └ Environment Variables
- |   └ Health Checks per Service
- |
- |
- | └ 10.3 Dockerfile Optimization
- |   └ Multi-stage Builds
- |   └ Layer Caching Strategy
- |   └ Image Size Optimization
- |   └ Security Best Practices
- |
- |
- | └ 10.4 Kubernetes Deployment (Ready)
- |   └ Deployment Manifests
- |   └ Service Definitions
- |   └ ConfigMaps & Secrets
- |   └ Horizontal Pod Autoscaling
- |   └ Rolling Updates Strategy
- |
- |
- └ 10.5 Infrastructure as Code
  - └ docker-compose.yml
  - └ Kubernetes YAML manifests
  - └ Environment Configuration
  - └ Deployment Scripts

---

## 11 RESULTADOS & MÉTRICAS DE SUCESSO (4 min)

### 11.1 Model Performance Results

- |   | Accuracy: 92.3% (test set)
- |   | F1-Score: 0.91 (balanced)
- |   | Precision: 0.93 (high confidence)
- |   | Recall: 0.90 (comprehensive)
- |   \ ROC-AUC: 0.96 (excellent discrimination)

### 11.2 System Performance Results

- |   | Latency Metrics
  - |   | P50: 42ms ⚡
  - |   | P95: 87ms ⚡
  - |   | P99: 145ms
  - |   \ Max: 250ms
- |   | Throughput: 1000+ req/min
- |   | Error Rate: < 0.1%
- |   | Uptime: 99.9%
- |   \ Cache Hit Rate: 78%

### 11.3 BERT vs GPT Comparison ★

- |   | Performance Comparison
  - |   | BERT: 45ms avg latency
  - |   | GPT: 1200ms avg latency
  - |   | Accuracy: BERT 92.3% vs GPT 93.1%
  - |   \ Agreement: 94.7% (Cohen's Kappa: 0.92)

- |   └ Cost Analysis
  - |   |   └ BERT: \$0/request (self-hosted)
  - |   |   └ GPT: ~\$0.0001/request
  - |   |   └ Savings: 99.9% using BERT
  - |   |
  - |   └ Recommendation Engine
    - |   |   └ BERT for 95% of cases
    - |   |   └ GPT for low-confidence (< 0.70)
    - |   |   └ Hybrid approach ROI
  - |   |
- |   └ 11.4 Business Impact Metrics
  - |   |   └ Total Predictions: 150K+
  - |   |   └ Active Users: Mock data
  - |   |   └ Feedback Collection Rate: 12%
  - |   |   └ Model Improvement Cycle: Weekly
- |   |
- |   └ 11.5 Technical Achievements
  - |   |   └ Test Coverage: 78%
  - |   |   └ Code Quality Score: A+
  - |   |   └ Documentation Pages: 15+
  - |   |   └ Time to Deploy: < 5 minutes

---

## 11 VANTAGENS COMPETITIVAS & DIFERENCIAIS (3 min)

- |   └ 12.1 Inovação Técnica
  - |   |   └  Comparação BERT vs GPT (ÚNICA!)
  - |   |   └ Endpoint dedicado
  - |   |   └ Análise de trade-offs
  - |   |   └ Recomendação inteligente

- | |
- | | └ LLM-as-Judge Framework
- | | | └ GPT valida BERT
- | | | └ Quality assurance automatizado
- | | | └ Continuous validation
- | |
- | | └ Explainability Built-in
- | | | └ LIME/SHAP integration
- | | | └ Word importance visualization
- | | | └ Confidence breakdown
- | |
- | | └ Continuous Learning Loop
  - | └ Feedback → Validation → Retrain
  - | └ Active learning ready
  - | └ Model versioning
- |
- | └ 12.2 Production-Ready desde o Dia 1
  - | └ Observability 360°
  - | └ Monitoring & Alerting
  - | └ Health Checks
  - | └ Circuit Breakers (Ready)
  - | └ Retry Logic
  - | └ Graceful Degradation
- |
- | └ 12.3 Escalabilidade & Performance
  - | └ Horizontal Scaling Ready
  - | └ Load Balancing (Nginx)
  - | └ Caching Strategy (Redis)

- |   └ Database Optimization
  - |   └ Async Processing
  - |
  - |   └ 12.4 Developer Experience
  - |    └ Deploy em 3 comandos
  - |    └ Documentação completa (15+ páginas)
  - |    └ Code Quality (78% coverage)
  - |    └ Clean Architecture
  - |    └ Onboarding < 5 minutos
  - |
  - |   └ 12.5 ROI & Custo-Benefício
  - |    └ 99.9% economia vs GPT-only
  - |    └ Latência 27x menor que GPT
  - |    └ Zero vendor lock-in
  - |    └ Privacidade garantida
  - |    └ Escalabilidade sem custos variáveis
- 

## **13 TROUBLESHOOTING & OPERAÇÃO (2 min)**

- |   └ 13.1 Common Issues & Solutions
  - |    └ Model Loading Failures
  - |    └ API Timeouts
  - |    └ Database Connection Issues
  - |    └ Cache Miss Spikes
  - |    └ High Latency Debugging
- |
- |   └ 13.2 Incident Response Playbook
  - |    └ Alert Investigation Workflow
  - |    └ Rollback Procedures

- |   └ Emergency Contacts
  - |   └ Post-mortem Template
  - |
  - |   └ 13.3 Performance Debugging
  - |   |   └ Profiling Tools (cProfile)
  - |   |   └ Bottleneck Identification
  - |   |   └ Query Optimization
  - |   |   └ Resource Tuning
  - |
  - |   └ 13.4 Maintenance & Updates
  - |   |   └ Model Retraining Schedule
  - |   |   └ Dependency Updates
  - |   |   └ Security Patches
  - |   |   └ Backup & Recovery
- 

## **14 ROADMAP & PRÓXIMOS PASSOS (2 min)**

- |   └ 14.1 Short-term (30 dias)
  - |   |   └ Cloud Deployment (AWS/GCP)
  - |   |   └ Authentication & Authorization
  - |   |   └ API Rate Limiting Enhancement
  - |   |   └ Mobile-friendly API
- |
- |   └ 14.2 Mid-term (60 dias)
  - |   |   └ Multi-language Support
  - |   |   └ Advanced Ensemble Models
  - |   |   └ Real-time Retraining Pipeline
  - |   |   └ Business Intelligence Layer
- |

- |— 14.3 Long-term (90+ dias)
    - |— Multi-tenant Architecture
    - |— Kubernetes Production Deployment
    - |— Advanced A/B Testing
    - |— Chatbot Integration
  - |— 14.4 Research & Innovation
    - |— Aspect-Based Sentiment Analysis
    - |— Multi-modal Analysis (text + images)
    - |— Emotion Detection
    - |— Sarcasm Detection
- 

## **15 CONCLUSÃO & CALL-TO-ACTION (2 min)**

- |— 15.1 Recapitulação
  - |— Problema → Solução → Resultado
  - |— Diferenciais Técnicos
  - |— Impacto de Negócio
- |— 15.2 Key Takeaways
  - |— Production-Ready desde o início
  - |— BERT vs GPT: melhor de ambos
  - |— Observability 360°
  - |— ROI comprovado
- |— 15.3 Demo Live (opcional)
  - |— Análise em tempo real
  - |— Dashboard Grafana
  - |— Comparação BERT vs GPT

|

└ 15.4 Q&A + Discussão

  └ Perguntas abertas

  └ Deep dive técnico

└ Próximos passos

---

 **RESUMO DE TEMPO SUGERIDO**

SEÇÃO	TEMPO (min)
1. Abertura & Contexto	3
2. Data Pipeline	4
3. Model Training	5
4. Model Validation & EVAL	4
5. Model Improvement	3
6. Backend API	4
7. Observability & Telemetry	4
8. Monitoring & Alerting	4
9. Testing Infrastructure	3
10. Containerization	3
11. Resultados & Métricas	4
12. Vantagens Competitivas	3
13. Troubleshooting	2
14. Roadmap	2
15. Conclusão & Q&A	2
TOTAL	50 min

---

## VERSÃO CONDENSADA (30 min)

Se você precisar de uma versão mais curta:

- Abertura (2 min)
  - Data + Training + EVAL (8 min) - Combinar 2, 3, 4
  - Backend + API (5 min)
  - Observability + Monitoring (6 min) - Combinar 7, 8
  - Testing + Docker (4 min) - Combinar 9, 10
  - Resultados + Vantagens (4 min) - Combinar 11, 12
  - Conclusão (1 min)
- 

## GLOSSÁRIO DE TERMOS TÉCNICOS USADOS

-  EVAL Framework: Sistema completo de avaliação de modelos
  -  Telemetry: Coleta automática de métricas e logs
  -  Troubleshooting: Diagnóstico e solução de problemas
  -  Observability: Capacidade de entender o sistema interno
  -  Alerting: Sistema de notificações automáticas
  -  Drift Detection: Identificação de mudanças em dados
  -  Continuous Learning: Aprendizado contínuo do modelo
  -  Fine-tuning: Ajuste fino de modelos pré-treinados
  -  LLM-as-Judge: LLM como validador de qualidade
  -  Explainability: Capacidade de explicar decisões
- 

**Esta agenda cobre 100% do projeto de forma profissional e técnica!** 