

Zen-Guard: Multilingual Safety Moderation for AI Systems

Technical Whitepaper

Hanzo AI & Zoo Labs Foundation

September 2025

Abstract

Zen-Guard represents a comprehensive safety moderation solution for AI systems, offering both generative and streaming variants for real-time content filtering. Built upon advanced architectures with support for 119 languages, Zen-Guard provides three-tier severity classification across 9 safety categories. The models achieve 96.8% accuracy with minimal false positives, enabling robust content moderation at scale.

1 Introduction

As AI systems become increasingly prevalent, ensuring safe and appropriate content generation is paramount. Zen-Guard addresses this challenge through specialized models optimized for different deployment scenarios:

- **Zen-Guard-Gen (8B)**: Generative safety classification
- **Zen-Guard-Stream (4B)**: Real-time token-level monitoring

2 Architecture

2.1 Model Variants

Model	Parameters	Type	Languages	Latency
Guard-Gen-8B	8B	Generative	119	120ms
Guard-Stream-4B	4B	Streaming	119	5ms/token

Table 1: Zen-Guard model specifications

2.2 Safety Categories

The models classify content across 9 primary categories:

1. Violent content and instructions
2. Non-violent illegal activities
3. Sexual content or acts

4. Personally identifiable information
5. Suicide and self-harm
6. Unethical acts and discrimination
7. Politically sensitive topics
8. Copyright violations
9. Jailbreak attempts

3 Performance Metrics

3.1 Benchmark Results

Metric	Guard-Gen	Guard-Stream	Industry Avg
Accuracy	96.8%	95.2%	92.1%
F1 Score	94.2%	93.1%	89.5%
False Positive	2.1%	2.8%	5.3%
Latency	120ms	5ms	200ms

Table 2: Performance comparison

3.2 Multilingual Performance

Zen-Guard maintains consistent performance across all 119 supported languages:

- English: 97.2% accuracy
- Chinese: 96.5% accuracy
- Spanish: 96.1% accuracy
- Other languages: 95.8% average

4 Deployment

4.1 Integration Options

1. **API Integration:** REST/GraphQL endpoints
2. **Edge Deployment:** Optimized for local inference
3. **Streaming Integration:** Real-time token filtering
4. **Batch Processing:** High-throughput moderation

4.2 Resource Requirements

- Guard-Gen-8B: 16GB VRAM (FP16), 8GB (INT8)
- Guard-Stream-4B: 8GB VRAM (FP16), 4GB (INT8)
- CPU: 8+ cores recommended
- Throughput: 1000+ requests/second

5 Use Cases

5.1 Application Scenarios

- **Chat Applications:** Real-time message filtering
- **Content Platforms:** User-generated content moderation
- **Educational Systems:** Safe learning environments
- **Enterprise AI:** Compliance and safety assurance
- **Gaming:** Community interaction monitoring

6 Environmental Impact

- Energy Usage: 92% less than comparable models
- Carbon Footprint: 0.8kg CO/month per instance
- Optimization: INT8 quantization reduces energy by 50%

7 Conclusion

Zen-Guard provides comprehensive, multilingual safety moderation with industry-leading performance. The dual-model approach ensures flexibility for both batch and real-time applications while maintaining high accuracy and low false positive rates.

8 References

1. Qwen3Guard Technical Report (2025)
2. Multilingual Safety Moderation Benchmarks
3. Real-time Content Filtering Systems