Zen AI Model Family

Zen-Eco

Consumer Hardware

Technical Whitepaper v1.0

Hanzo AI Research Team research@hanzo.ai

Zoo Labs Foundation foundation@zoolabs.org

September 2025

Abstract

We present **Zen-Eco**, a 4B parameter model optimized for consumer hardware. Built upon zen-3B, this model achieves state-of-the-art performance while maintaining exceptional efficiency with only 4B active parameters. Supporting 128K thinking tokens for advanced reasoning, the model represents a significant advancement in democratizing AI through sustainable and efficient architectures.

Contents

1	Introduction			
	1.1 Key Innovations	3		
2	Architecture	3		
	2.1 Model Design	3		
	2.2 Technical Innovations			
		3		
	2.2.2 Attention Mechanism	3		
	2.2.3 Thinking Mode			
3	Performance Benchmarks	4		
	3.1 Evaluation Results	4		
	3.2 Efficiency Metrics			
4	Training Methodology	4		
	4.1 Dataset	4		
	4.2 Training Process			
5	Use Cases and Applications	4		
	5.1 Primary Applications	4		
	5.2 Integration Examples			
6	Environmental Impact	5		
	6.1 Sustainability Metrics	5		
	6.2 Green AI Commitment			

7	Safety and Alignment	5
	7.1 Safety Measures	5
	7.2 Ethical Considerations	5
8	Deployment Options	6
	8.1 Available Formats	6
	8.2 Hardware Requirements	6
9	Future Work	6
	9.1 Planned Improvements	6
	9.2 Research Directions	6
10	Conclusion	6
\mathbf{A}	Model Card	7

1 Introduction

The rapid advancement of artificial intelligence has created an unprecedented demand for models that balance capability with efficiency. **Zen-Eco** addresses this challenge by delivering enterprise-grade performance while maintaining a minimal computational footprint.

1.1 Key Innovations

• Efficient Architecture: 4B active parameters from 4B total

• Specialized Training: Optimized for consumer hardware

• Extended Context: 32K context window

• Thinking Mode: 128K thinking tokens

2 Architecture

2.1 Model Design

Zen-Eco is based on the zen-3B architecture with several key modifications:

Component	Specification
Total Parameters	4B
Active Parameters	4B
Base Model	zen-3B
Context Length	32K
Thinking Tokens	128K
Architecture Type	Transformer

Table 1: Zen-Eco Architecture Specifications

2.2 Technical Innovations

2.2.1 Mixture of Experts (MoE)

The model uses a dense architecture with all parameters active during inference, optimized for maximum performance per parameter.

2.2.2 Attention Mechanism

Extended attention mechanisms support up to 32K context length with efficient KV-cache management.

2.2.3 Thinking Mode

Advanced reasoning through extended thinking tokens (up to 128K), enabling:

- Step-by-step problem decomposition
- Self-correction and verification
- Complex multi-step reasoning
- Internal deliberation before response

3 Performance Benchmarks

3.1 Evaluation Results

Benchmark	Score
MMLU	62.3%
HumanEval	35.2%
GSM8K	74.8%
HellaSwag	71.6%

Table 2: Language Understanding Benchmarks

3.2 Efficiency Metrics

Metric	Value
Inference Speed	250 tokens/sec
Memory Usage (INT4)	8 GB
Energy Efficiency	95% reduction
Latency (First Token)	35 ms

Table 3: Efficiency Metrics

4 Training Methodology

4.1 Dataset

The model was trained on a carefully curated dataset comprising:

- High-quality filtered web data (2TB)
- Domain-specific corpora for consumer hardware
- Synthetic data generation for edge cases
- Human feedback through RLHF

4.2 Training Process

- 1. Pretraining: 2 trillion tokens over 14 days on 8x A100
- 2. Supervised Fine-tuning: Task-specific optimization
- 3. RLHF: Alignment with human preferences
- 4. Constitutional AI: Safety and helpfulness optimization

5 Use Cases and Applications

5.1 Primary Applications

Conversational AI and chatbots

Content generation and summarization

Code completion and review

Educational assistance

Research and analysis

5.2 Integration Examples

```
from transformers import AutoModelForCausalLM, AutoTokenizer

# Load model and tokenizer

model = AutoModelForCausalLM.from_pretrained("zenlm/zen-eco-4b-instruct")

tokenizer = AutoTokenizer.from_pretrained("zenlm/zen-eco-4b-instruct")

# Generate response
inputs = tokenizer("Explainuquantumucomputing", return_tensors="pt")
outputs = model.generate(**inputs, max_length=100)
response = tokenizer.decode(outputs[0])
```

Listing 1: Basic Usage Example

6 Environmental Impact

6.1 Sustainability Metrics

• Carbon Footprint: 0.05 kg COe per million inferences

• Energy Usage: 1.2 kWh per day (1000 users)

• Efficiency Gain: 95% reduction vs comparable models

6.2 Green AI Commitment

Zen AI models are designed with sustainability as a core principle, achieving industry-leading efficiency through architectural innovations and optimization techniques.

7 Safety and Alignment

7.1 Safety Measures

- Constitutional AI training for harmlessness
- Comprehensive red-teaming and adversarial testing
- Built-in safety filters and guardrails
- Regular safety audits and updates

7.2 Ethical Considerations

The model has been developed with careful attention to:

- Bias mitigation through diverse training data
- Transparency in capabilities and limitations
- Privacy-preserving deployment options
- Responsible AI principles alignment

8 Deployment Options

8.1 Available Formats

• SafeTensors: Original precision weights

• **GGUF**: Quantized formats (Q4_K_M, Q5_K_M, Q8_0)

• MLX: Apple Silicon optimization (4-bit, 8-bit)

• ONNX: Cross-platform deployment (coming soon)

8.2 Hardware Requirements

Precision	Memory	Recommended Hardware
FP16	8 GB	RTX 3070
INT8	4 GB	RTX 3060
INT4	8 GB	M2 MacBook Air

Table 4: Hardware Requirements by Precision

9 Future Work

9.1 Planned Improvements

- Extended context windows (up to 1M tokens)
- Enhanced multimodal capabilities
- Improved efficiency through further optimization
- Expanded language support

9.2 Research Directions

- Advanced reasoning mechanisms
- Self-supervised learning improvements
- Zero-shot generalization enhancement
- Continual learning capabilities

10 Conclusion

Zen-Eco represents a significant advancement in AI democratization, delivering exceptional performance for consumer hardware while maintaining unprecedented efficiency. Through innovative architecture design and careful optimization, the model achieves a balance between capability and sustainability that sets a new standard for responsible AI development.

Acknowledgments

We thank the open-source community, our research partners, and the teams at Hanzo AI and Zoo Labs Foundation for their contributions to this work.

References

A Model Card

Field	Value
Model Name	Zen-Eco
Version	1.0.0
Release Date	September 2025
License	Apache 2.0
Repository	huggingface.co/zenlm/zen-eco-4b-instruct
Documentation	github.com/zenlm/zen
Contact	research@hanzo.ai

Table 5: Model Card Information