# AI COST MANAGEMENT: A STRATEGIC APPROACH

Essential Strategies for Al Provider Cost Optimization Managing OpenAl, Anthropic & Multi-Model Environments

#### **EXECUTIVE SUMMARY**

Al costs are spiraling out of control for most organizations, with companies spending 5x more than necessary on Al inference and training. This guide provides actionable strategies to optimize Al spending while maintaining quality and performance across multiple Al providers.

## **KEY RESEARCH FINDINGS:**

⢠75% average cost reduction through intelligent model routing
⢠40% token savings via prompt optimization
⢠90% reduction in budget overruns with governance frameworks
⢠Support for 5+ AI providers including OpenAI, Anthropic, Cohere

## TABLE OF CONTENTS

Chapter 1: Al Cost Landscape Analysis	Pages 3-6
Chapter 2: Model Selection Framework	Pages 7-12
Chapter 3: Token Optimization Strategies	Pages 13-18
Chapter 4: Multi-Model Routing	Pages 19-24
Chapter 5: Al Governance & Budget Control	Pages 25-29
Chapter 6: Case Studies & ROI Analysis	Pages 30-32

## **OPTIMIZATION TECHNIQUES**

- Intelligent Model Selection
   Choose optimal model for each task based on cost/performance ratio
- Dynamic Prompt Optimization Reduce token consumption through prompt engineering
- Multi-Provider Routing Route requests to most cost-effective provider in real-time
- Context Window Management
   Optimize context usage to minimize token costs
- Caching & Response Reuse Implement intelligent caching to avoid redundant API calls

#### PROVIDER COMPARISON

OpenALCDT 4: Dromium quality, higher cost per taken