✦ Fireworks AI

# Unlocking Multimodal AI for Vision & Text Intelligence

Actionable Insights, Faster Decision-Making, and Scalable Automation

---

## Executive Summary                                      01

---

Enterprises are overwhelmed with unstructured data across product images and documents. Most AI tools handle only text, leaving critical insights trapped, slowing workflows, and driving up costs. By 2025, 80% of the global datasphere is projected to be unstructured ([IDC, Data Age 2025](#)), highlighting the growing risk of blind spots for enterprises without domain-tuned AI.

Fireworks AI provides a production-ready platform for scalable, fine-tuned multimodal AI intelligence, enabling enterprises to extract structured, actionable insights and automate workflows reliably at scale. Key capabilities include:

- **Vision & Text Intelligence:** Fine-tuned vision-language models (VLMs) for product classification, document parsing, and structured data extraction at scale.
- **Domain Adaptation & Fine-Tuning:** Fine-tune vision-language models aligned to enterprise taxonomies, document schemas, and internal knowledge for higher accuracy, efficiency, and automation.

Enterprises adopting Fireworks AI can expect to:

- Achieve fast, context-aware understanding across images and documents.
- Scale low-latency processing for production workloads.
- Reduce operational costs, increase automation, and drive measurable business impact.

---

### The Case for AI in Software Development

**MACRO TREND** — Enterprises are drowning in unstructured data. Analysts spend hours manually tagging images and parsing forms, while decisions slow down and costs rise. Without multi-modal AI, valuable insights remain locked in raw content.

**CHALLENGE** — Most enterprises lack domain-tuned AI that can handle images and documents with domain-specific accuracy. As a result, critical insights remain locked, slowing workflows and increasing costs.

**SOLUTION** — Fireworks provides enterprise-grade multimodal AI that can extract insights from text and images out-of-the-box. Turn unstructured data into formatted text that can be systematically processed at scale.

# Fireworks AI: The Platform for Multimedia Intelligence

Fireworks AI provides a complete, enterprise-ready platform for processing images, documents, and product data in real time.

| Feature | Enterprise Value | How Fireworks Enables It | Metrics/Benchmarks |
|---|---|---|---|
| **Model choice & customization** | Tailored AI for domain-specific data | Choose between different VLMs and speech models or use fine-tuned models | Accuracy on domain-specific tasks |
| **Enterprise-ready infrastructure** | Low-latency, high-throughput inference, reliable scale for production workloads | Optimized, proprietary serving stack, GPU auto-scaling, and global infrastructure | Sub-2s latency, millions of queries/day |
| **Cost Efficiency** | Lower TCO for large-scale AI deployments | Optimized serving architecture that reduces GPU/infra waste | Cost-optimized deployment |
| **Control & Flexibility** | Enterprise governance and compliance controls | BYOC deployment for vision models, FireOptimizer, audit and RBAC controls | Enterprise-grade compliance and traceability for regulated workloads |
| **AI Orchestration/ Gateway** | Centralized model management | Unified API, model catalog, monitoring | Enabled consistent operations across thousands of endpoints |
| **Model Lifecycle Management** | Continuous performance and cost optimization | Build → Tune → Scale cycles | Consistent throughput, reduced downtime, cost-optimized deployment |

This combination of **speed, precision, cost efficiency, and enterprise control** addresses the key challenges enterprises face when extracting insights, and automating workflows across images and documents.

# Patterns & Practices for Multimedia AI

## People & Process

- Build expertise to fine-tune models for domain-specific image and documents.
- Integrate AI-assisted workflows into product and operational pipelines.
- Collaborate between AI/ML teams and product/platform teams.

## Technology Implementation

- Deploy VLMs models on scalable infrastructure.
- Use Fireworks AI Gateway to centralize orchestration, monitoring, and model upgrades.
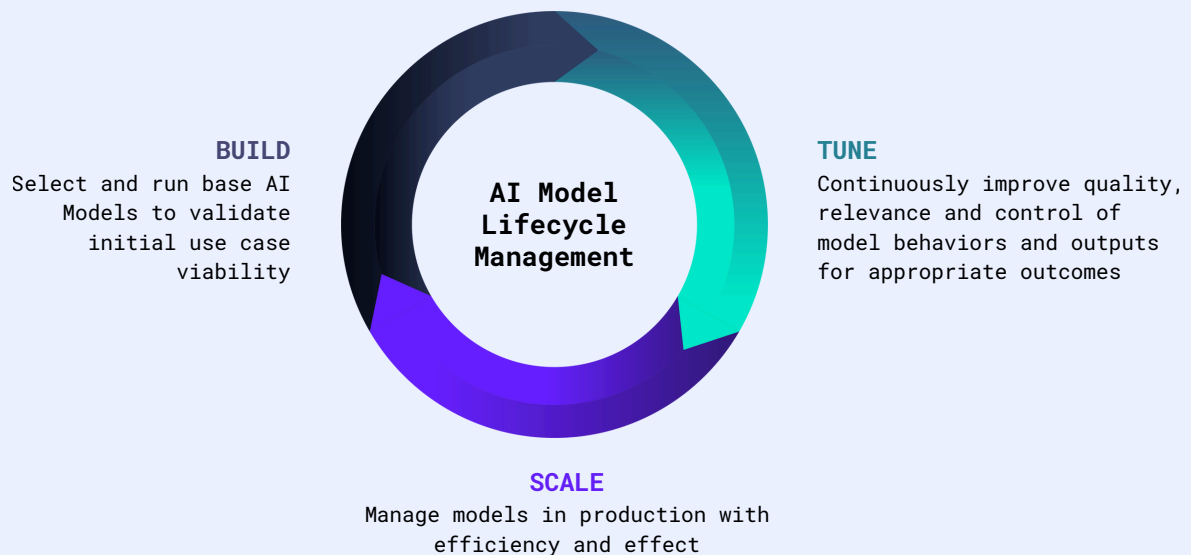- Continuously evaluate performance, accuracy, and cost efficiency.

## Risk Management

- Ensure compliance and auditability with BYOC, RBAC, and logging.
- Protect proprietary data and model outputs.
- Minimize vendor lock-in and operational risk.

## Operational Patterns

- **AI Gateway:** Unified API for inference, monitoring, and scaling.
- **Model Lifecycle Management:** Build → Tune → Scale selection, fine-tuning, continuous deployment, cost optimization.
- **Product-Model Co-Design:** Embed multimodal AI capabilities into workflows and products for measurable ROI.

## Model Lifecycle Management

**BUILD**
Select and run base AI Models to validate initial use case viability

**AI Model Lifecycle Management**

**TUNE**
Continuously improve quality, relevance and control of model behaviors and outputs for appropriate outcomes

**SCALE**
Manage models in production with efficiency and effect

Multimodal AI enables enterprises to:

- **Accelerate decision-making** and operational efficiency.
- **Reduce manual effort** in tagging, classification, and extraction.
- **Ensure compliance and auditability** in sensitive workflows.
- **Scale insights** across millions of SKUs and documents.
- **Drive smarter automation** and product experiences with domain-specific intelligence.

Deploying scalable, low-latency multimodal AI allows organizations to unlock faster, more accurate, and cost-effective insights across text and images, enabling reliable automation and informed business decisions at enterprise scale.

## Getting Started

**1**

**Identify high-value workflows**

Start with targeted use cases such as Product catalog enrichment and document parsing.

**2**

**Run a pilot**

Deploy Fireworks AI with recommended models to validate performance, latency, and cost efficiency.

**3**

**Scale and customize**

Fine-tune models using enterprise-specific data and feedback loops.

**4**

**Implement operational patterns**

Adopt AI Gateway, Model Lifecycle Management, and product-model co-design.

# Conclusion

Enterprises can unlock real-time, structured insights from unstructured data at scale with Fireworks AI. By combining domain-tuned vision-language models with enterprise-grade infrastructure and lifecycle management, Fireworks AI transforms image and document data into actionable intelligence faster, more accurately, and cost-effectively.

## Next Steps

Fireworks AI helps enterprises move from experimentation to production with confidence. To get started:

| | | |
|---|---|---|
| | **Assess workflows** | Identify high-value tasks involving images and documents. |
| | **Run a pilot** | Deploy Fireworks AI with fine-tuned multimodal for reliable extraction. |
| | **Fine-tune and scale** | Customize models for domain-specific datasets and scale across production workloads. |
| | **Adopt operational patterns** | Implement AI Gateway, Model Lifecycle Management, and product-model co-design. For a broader perspective on enterprise AI architecture and operational maturity, see our enterprise page. |

## Fireworks AI

**Contact Fireworks AI today** to schedule a pilot, explore model customization, and unlock scalable, high-performance code intelligence for your development teams.