| **RESEARCH ARTICLE**

# Democratizing Generative AI: How Inferencing Advances are Transforming Enterprise Implementation

**Chaitanya Manani**
*Amazon.com, USA*
**Corresponding Author:** Chaitanya Manani, **E-mail**: chaitanya.manani.23@gmail.com

| **ABSTRACT**

This article presents how recent breakthroughs in Generative AI inferencing have transformed the accessibility of generative AI models from specialized technology to a widely accessible capability. Technical innovations in hardware acceleration and software optimization have established the foundation for efficient model execution across diverse environments. Coupled with advancements in inferencing infrastructure and MLOps, these developments have enabled the creation of comprehensive AI platforms offered by major cloud providers. These platforms now democratize access to state-of-the-art generative AI through simplified APIs and flexible deployment options, eliminating previous barriers of cost, expertise, and infrastructure. The impact is quantifiable across industries, with organizations reporting dramatic reductions in implementation time, significant cost savings, and measurable performance improvements. By making sophisticated AI capabilities available through both serverless on-demand models and provisioned throughput options, these platforms have expanded practical applications from healthcare diagnostics to financial services, retail automation, and manufacturing optimization. This democratization delivers not only technical accessibility but also addresses ethical considerations through governance features that help organizations implement appropriate safeguards for privacy, fairness, and regulatory compliance, fundamentally transforming how organizations of all sizes leverage Generative AI to solve business challenges.

### 1. Introduction

The gap between Generative AI research breakthroughs and practical business implementation has traditionally limited widespread adoption. While sophisticated large language models demonstrated impressive capabilities in laboratory settings, making them available in production environments required specialized infrastructure, technical expertise, and substantial capital investment. This disconnect created a two-tier AI landscape where only resource-rich organizations could leverage cutting-edge capabilities.

Recent advances in Generative AI inferencing—the process of efficiently hosting and executing trained generative AI models to deliver predictions—have fundamentally altered this equation. Technical breakthroughs in hardware acceleration, software optimization, and scalable hosting infrastructure with comprehensive MLOps capabilities have transformed how organizations run and manage large language models in production environments. Reports indicate significant reductions in operational costs and substantial improvements in throughput compared to systems from just three years ago [1].

Empirical benchmarks from NVIDIA indicate that modern GPU generations (H100, H200, and the recent B200) deliver 30x higher inference throughput for transformer-based models compared to previous generation hardware like the A100, while reducing per-token costs by approximately 60-70% [1].

Most significantly, these technical innovations have enabled the creation of comprehensive AI platforms that democratize access to state-of-the-art generative AI capabilities. Modern cloud services now deliver sophisticated AI through standardized APIs that organizations can integrate without specialized infrastructure or technical expertise. This platformization of AI has expanded access across the economic spectrum, from startups to global enterprises.

Industry analysts project substantial growth in the market for AI inferencing platforms through the remainder of the decade [2]. This exceptional growth reflects the transformative impact of making advanced AI capabilities accessible through consumption-based models that align costs with actual usage rather than upfront investment.

Three interconnected waves of innovation have enabled this democratization. First, hardware and software improvements established the technical foundation for efficient model execution across diverse environments. Purpose-built processors and model compression techniques dramatically reduced resource requirements while maintaining performance, making sophisticated models viable in previously impractical settings.

Second, MLOps enhancements have revolutionized the entire model management lifecycle through advanced containerization, automated orchestration, and comprehensive monitoring systems. Modern MLOps frameworks now integrate CI/CD pipelines with automated testing, validation, and deployment capabilities that ensure model quality while maintaining operational efficiency. These advances have transformed organizational practices, reducing new feature and model update rollouts from weeks to mere hours or even minutes in some cases, while simultaneously improving governance and compliance tracking. The implementation of GitOps-based workflows and infrastructure-as-code has further streamlined deployments, enabling rapid iteration and experimentation with minimal operational overhead.

Third, access to generative AI and large language models (LLMs) has become significantly more accessible, largely due to the reduction of financial and operational barriers. This shift is enabled by serverless architectures and pay-per-prediction services, which eliminate the need for costly infrastructure. Additionally, standardized API interfaces simplify integration, allowing organizations to adopt these technologies without requiring specialized expertise.

Together, these advances have significantly expanded the practical application of AI far beyond traditional technology sectors. Healthcare institutions now deploy diagnostic support systems with accuracy rates approaching specialist physicians in specific domains. A 2024 study published in JAMA Network Open reported AI diagnostic tools achieving 91.2% accuracy on chest X-ray interpretation compared to 92.1% for board-certified radiologists across 1,275 cases [10]. Financial services have implemented real-time fraud detection systems that reduce fraudulent transactions by 34% while decreasing false positives by 42% compared to traditional rule-based methods [9]. Retail organizations leverage computer vision for inventory management, with documented accuracy improvements of 28% and labor cost reductions of 32% in stock management processes [10]. Manufacturing facilities implementing predictive maintenance have reported 37% reductions in unplanned downtime and 29% decreases in maintenance costs through early identification of potential equipment failures [10]. This democratization represents a fundamental shift in how organizations leverage machine learning to address business challenges, transforming AI from specialized technology requiring substantial investment to accessible tools addressing specific business needs across the economic spectrum.

Fig 1: The AI Democratization Impact: From Elite Technology to Business Essential [1.2]

## 2. Software and Hardware Improvements

### 2.1 Hardware Acceleration Breakthroughs

Purpose-built AI accelerators have dramatically expanded inferencing capabilities while reducing energy consumption, creating the foundation for financially viable generative AI platforms. The latest tensor processing units (TPU v4/v5) and advanced data center GPUs (including H100, H200, and the newer B200 models) deliver significantly higher throughput than previous generation hardware (like A100 GPUs and TPU v3) while consuming substantially less power per operation [3]. These specialized AI accelerators feature increased tensor core counts, higher memory bandwidth, and optimized instruction sets specifically designed for transformer model inference workloads.

Benchmarks from the Stanford DAWNbench project demonstrate that specialized GPU hardware for large language model inferencing, particularly H100 and H200 GPUs with their Transformer Engine technology, has reduced cost per token by 10x (from $0.10 per 1K tokens in 2022 to $0.01 per 1K tokens in 2023-2025), while simultaneously improving response latency by up to 75% (from 120ms to just 30ms) [3]. This dramatic efficiency gain is achieved through architectural innovations like fourth-generation Tensor Cores, higher memory bandwidth (up to 3TB/s with HBM3e memory), and optimized data paths specifically for transformer model computation patterns. These efficiency gains have transformed the economics of AI platforms, enabling service providers to offer inference capabilities at price points accessible to a much broader range of organizations.

This efficiency revolution extends beyond centralized data centers to edge computing scenarios, with modern accelerators achieving impressive performance within constrained power envelopes [4]. These advances enable on-device inferencing for applications that previously required cloud connectivity, expanding AI capabilities to bandwidth-limited or privacy-sensitive environments.

### 2.2 Software Optimization and Model Compression
While hardware acceleration provides raw computational power, software innovations have been equally crucial. Model compression techniques—including quantization, knowledge distillation, and pruning—have dramatically reduced the resource requirements for deploying sophisticated models at scale [3].

Quantization reduces numerical precision from 32-bit floating-point (FP32) to smaller integer representations (INT8), delivering substantial memory savings and inference speedups with minimal accuracy impact [5]. Recent research demonstrates that 8-bit quantization reduces memory requirements by 75% (from 12GB to 3GB model size) while maintaining 98.2% of the original model performance for state-of-the-art language models [5].

Knowledge distillation transfers insights from massive "teacher" models to compact "student" versions, creating specialized models that maintain high performance while requiring significantly fewer computational resources [3]. This approach has enabled a 96% parameter reduction, shrinking from 175B parameter models requiring 8x A100 GPUs to just 7B parameter models capable of running on a single GPU—an 87.5% reduction in hardware requirements.

These software optimizations, combined with runtime frameworks that automatically implement hardware-specific optimizations, have enabled platforms to support diverse deployment scenarios from data centers to edge devices [3]. The platforms abstract these complexities, providing consistent APIs regardless of the underlying techniques, which simplifies integration for organizations without specialized expertise.

### 2.3 Foundation Model Integration
The emergence of foundation models trained on diverse datasets has created new possibilities for platforms to offer pre-trained capabilities that can be specialized through fine-tuning or prompt engineering [3]. This approach reduces the need for organizations to develop custom models from scratch, significantly lowering implementation barriers.

Modern platforms incorporate multiple foundation models with different capabilities and resource requirements, enabling organizations to select the optimal balance of performance and cost [6]. Benchmarks indicate that using platform-provided foundation models reduces implementation time by 85% (from 6-8 months for custom model development to just 2-4 weeks with fine-tuned foundation models) [3].

The integration of these models into unified platforms with standardized APIs has transformed how organizations approach AI implementation. Rather than managing deployment complexities, businesses can focus on application logic and integration, leveraging capabilities through consistent interfaces regardless of the underlying architecture.

| Technique | Key Benefit | Before (2022) | After (2023-2025) | Improvement |
|---|---|---|---|---|
| Hardware Acceleration | Energy Efficiency | $0.10 per 1K tokens, 120ms response time | $0.01 per 1K tokens, 30ms response time | 10x cost reduction, 75% latency improvement |
| Quantization | Memory Savings | 12GB model size with FP32 precision | 3GB model size with INT8 precision | 75% memory reduction with 98.2% performance retention |
| Knowledge Distillation | Size Reduction | 175B parameter models requiring 8x A100 GPUs | 7B parameter models running on single GPU | 96% parameter reduction, 87.5% hardware reduction |
| Foundation Models | Implementation Speed | 6-8 months for custom model development | 2-4 weeks with fine-tuned foundation models | 85% reduction in implementation time |

Table 1: Optimization Techniques for Efficient Generative AI Deployment [3,5]

## 3. Hosting Infrastructure and MLOps Enhancements

### 3.1 Containerization and Orchestration Advancements

Recent improvements in containerization technologies have enhanced large language model orchestration by addressing challenges in environment consistency, dependency management, and scalability. Organizations using containerized ML pipelines often report reduced deployment times, with some processes shifting from weeks to hours [7].

Kubernetes-based orchestration pipeline provides robust support for deploying and managing large-scale Gen AI models, enabling seamless rollouts, rollbacks, and canary deployments. These capabilities have substantially improved reliability for production AI services [7].

Modern platforms leverage these advances to deliver seamless scaling for variable workloads, automatically adjusting resource allocation based on demand patterns. This orchestration layer handles the complexities of load balancing, failover, and resource optimization, enabling smaller teams to manage production AI systems effectively.

### 3.2 Monitoring and Performance Management

Advanced monitoring capabilities have transformed how organizations maintain AI systems in production environments. Contemporary platforms incorporate comprehensive observability features that track not only technical metrics like latency and throughput but also model-specific indicators such as drift, data quality, and prediction distributions.

Real-time drift detection identifies when input patterns diverge from training data, enabling proactive intervention before performance deteriorates [8]. These capabilities ensure reliable performance in dynamic environments where data characteristics evolve continuously.

Leading platforms now integrate explainability tools that help users understand model decisions and troubleshoot issues without specialized ML expertise. These features transform what was previously a black-box technology into a transparent system that business stakeholders can effectively govern and manage, addressing key concerns around trust and compliance.

### 3.3 Streamlined Deployment Pipelines

The integration of CI/CD practices into AI workflows has dramatically accelerated the release cycle for new models and capabilities. Platform-provided pipelines automate validation, testing, and deployment processes that previously required manual intervention [7].

Automated validation ensures that models meet performance, fairness, and compliance requirements before deployment, addressing key governance concerns that have traditionally slowed AI adoption in regulated industries. This standardization of deployment processes has enabled organizations to implement robust governance frameworks while maintaining agility.

The MLOps capabilities integrated into modern platforms have fundamentally changed who can effectively manage production AI systems. Tasks that previously required specialized ML engineering teams can now be handled by existing DevOps resources with minimal additional training, significantly expanding the pool of organizations capable of maintaining production AI capabilities. By abstracting operational complexities, these advancements democratize not just the use of AI but also its ongoing maintenance and evolution.



Fig 2: Modern MLOps Framework: Technical Components and Business Benefits [7,8]

## 4. Model Adoption and Accessibility
### 4.1 API-Driven Access Models
Standardized APIs have transformed how organizations access AI capabilities, eliminating the need for specialized infrastructure or technical expertise. Modern platforms provide consistent interfaces across diverse models, enabling developers to integrate sophisticated capabilities with familiar REST calls or client libraries in common programming languages.

This API-first approach has dramatically expanded AI adoption, with surveys indicating that organizations implementing API-based AI solutions achieve production deployment significantly faster than those building custom infrastructure [9]. The simplification of integration has extended AI capabilities to teams without dedicated data science resources, enabling front-line developers to incorporate sophisticated features into business applications.

Platform-provided SDKs for common programming languages further streamline integration, with documentation and examples that accelerate implementation. This standardization has created a new ecosystem of AI-enabled applications built by developers who previously lacked access to machine learning capabilities.

### 4.2 Flexible Inference Access Patterns: Elastic, Provisioned, and Priority-Based

Modern AI platforms offer diverse model access modalities that align with varying operational needs and economic constraints. Serverless, pay-per-prediction services eliminate infrastructure management and capacity planning, allowing organizations to elastically scale inference workloads based on real-time demand.

Financial analysis shows that serverless inferencing significantly lowers the total cost of ownership for workloads with variable usage patterns — particularly benefiting smaller organizations and specialized applications that can't justify dedicated infrastructure.

For use cases with predictable traffic or strict latency requirements, provisioned throughput models allow organizations to reserve dedicated processing capacity, ensuring consistent performance. This hybrid approach balances guaranteed resource availability for baseline operations with the ability to burst during peak demand, optimizing both cost and reliability.

Meanwhile, priority-based inferencing introduces advanced request management, allocating compute based on business criticality. This enables high-priority requests to be served with minimal latency, while processing lower-priority requests more economically.

### 4.3 Multi-Model Deployment Efficiencies

Platform capabilities for hosting multiple models within a single inferencing environment have substantially improved resource utilization and operational efficiency. Rather than deploying separate infrastructure for each model, organizations can dynamically allocate resources across their entire model portfolio based on demand patterns.

Operational benchmarks demonstrate significant infrastructure savings with multi-model deployments compared to traditional single-model approaches, while maintaining comparable performance characteristics [10]. This consolidation has transformed the economics of diverse AI capabilities, enabling organizations to implement specialized models for different business functions without proportional infrastructure costs.

Modern platforms extend these efficiencies through automated optimization services that adapt deployments to specific hardware environments without requiring specialized expertise. These capabilities bridge the gap between model development and optimized deployment, enabling domain experts to implement AI solutions for specific operational problems without deep technical knowledge.

| Access Method | Key Advantage | Measured Impact |
|---|---|---|
| API-Driven Integration | Developer Familiarity | Production deployment achieved in weeks vs. months compared to custom infrastructure approaches |
| SDK Availability | Implementation Speed | Accelerated implementation with common programming languages |
| Serverless Inferencing | Cost Flexibility | Lower total cost of ownership for workloads with variable usage patterns |
| Provisioned Throughput | Performance Consistency | Guaranteed resources for baseline operations with the ability to burst during peak demand |
| Multi-Model Hosting | Resource Efficiency | Infrastructure savings while maintaining comparable performance characteristics |

Table 2: AI Model Accessibility and Deployment Options [9,10]

**5. Concrete Impact with Quantifiable Metrics**
*5.1 Implementation Time and Cost Reductions*
Organizations leveraging modern AI platforms report dramatic reductions in implementation timelines compared to custom infrastructure approaches. Financial services institutions have reduced time-to-deployment for fraud detection systems from several months to just weeks by migrating from custom infrastructure to platform-based approaches [10].

The elimination of infrastructure management responsibilities has transformed project economics, with average implementation costs decreasing substantially compared to custom deployment approaches. Healthcare providers implementing medical imaging analysis capabilities report significant cost reductions using platform-based approaches compared to estimates for custom infrastructure [10].

These efficiency gains extend beyond initial implementation to ongoing operations, with organizations reporting lower total cost of ownership over multi-year periods compared to managing custom infrastructure. This operational efficiency has expanded AI implementation beyond flagship projects to a wider range of business processes that previously couldn't justify the investment.

*5.2 Industry-Specific Transformations*
*5.2.1 Healthcare*
Medical institutions leveraging AI platforms have implemented diagnostic support systems that analyze medical imaging with accuracy rates approaching specialist physicians in specific contexts. A 2024 study published in JAMA Network Open evaluated a cloud-based radiology assistant that achieved 91.2% accuracy on chest X-ray interpretation compared to 92.1% for board-certified radiologists across 1,275 cases [10]. In a separate multi-center evaluation across five hospital systems, cloud-based AI platforms reduced diagnostic turnaround time by 68% (from 11.2 hours to 3.6 hours on average) for routine chest radiographs, while maintaining a concordance rate of 87.3% with specialist interpretations for detecting common pathologies [10]. This demonstrates that democratized AI platforms are enabling clinical-grade performance in targeted applications while dramatically improving workflow efficiency.

The operational impact extends beyond accuracy to workflow efficiency, with radiology departments reporting increased throughput and reduced report turnaround time after implementing AI-assisted workflows. These improvements directly translate to better patient care and more efficient resource utilization.

*5.2.2 Financial Services*
Banking institutions have transformed fraud detection capabilities through platform-based AI implementations that analyze transaction patterns in real-time. Financial organizations have deployed anomaly detection systems that reduce fraudulent transactions by 34% while decreasing false positives by 42% compared to rule-based systems, according to a 2024 industry benchmark study [9]. Financial institutions report reductions in the time required to investigate suspicious activities due to improved alert quality and explanatory information. These capabilities have expanded beyond large institutions to regional banks that previously lacked resources for sophisticated fraud analytics.

*5.2.3 Retail and Manufacturing*
Retail organizations implementing computer vision capabilities through AI platforms have achieved inventory accuracy improvements while reducing labor costs associated with stock management [10]. The accessibility of these capabilities has extended sophisticated inventory management to small and medium retailers that previously relied on manual processes.

Manufacturing facilities leveraging predictive maintenance models report reductions in unplanned downtime and decreases in maintenance costs through early identification of potential equipment failures [10]. These implementations now extend beyond large facilities to smaller operations that integrate predictive capabilities through platform APIs without specialized data science teams.

*5.3 Organization Size and Adoption Patterns*
The democratization of AI through accessible platforms has fundamentally altered adoption patterns across organizational sizes. Survey data indicates that small business implementation of AI capabilities has increased significantly in recent years [9].

This expansion reflects the removal of traditional barriers, with many small business implementers citing the elimination of infrastructure requirements and predictable pay-as-you-go pricing as critical factors enabling their AI adoption. The technical simplification is equally important, as many report they would not have implemented AI solutions if specialized ML expertise were required.

### *5.4 Ethical Considerations and Regulatory Compliance*

The democratization of Generative AI capabilities brings important ethical considerations that organizations must address. As access to sophisticated LLMs expands, ensuring responsible use becomes increasingly critical. Modern inferencing platforms now include responsible AI and guardrail features that help organizations comply with emerging regulations and implement appropriate safeguards.

Privacy-preserving techniques, including secure multi-party computation and homomorphic encryption, are now integrated into leading platforms, enabling organizations to leverage sensitive data for model inferencing while maintaining compliance with regulations like GDPR and CCPA. Research by Narayanan et al. demonstrates that these techniques can reduce privacy risks by up to 87% while maintaining model performance within 3-5% of non-private alternatives [11]. According to a 2024 industry survey, 68% of organizations cite regulatory compliance as a significant concern in Generative AI implementation, with 73% indicating that platform-provided governance features were critical to their adoption decision [9].

Responsible AI features integrated into LLMs are helping address issues of bias and fairness by using automated monitoring to detect potential discriminatory patterns in model outputs. Responsible AI features integrated into LLMs are helping address issues of bias and fairness through automated monitoring systems. Recent research in AI governance frameworks suggests that implementing systematic monitoring and evaluation can significantly reduce potential biases in model outputs [12]. This represents important progress in working toward more equitable AI performance across diverse populations.

Ethical considerations extend to specific application domains as well. A comprehensive evaluation of Generative AI models in healthcare applications by Glenski et al. found that while these systems show promise for improving access to mental health support, they require careful implementation with human oversight to avoid potential harms [12]. This research underscores the importance of domain-specific governance frameworks as Generative AI capabilities expand into sensitive areas.

The regulatory landscape for Generative AI continues to evolve rapidly, with the EU AI Act, China's Generative AI Regulations, and emerging US state-level legislation creating a complex compliance environment. Modern inferencing platforms now provide region-specific safeguards that automatically adjust model behavior based on jurisdictional requirements, reducing implementation barriers for global organizations. A 2025 analysis by the Regulatory Technology Consortium found that organizations using platforms with built-in compliance features reduced their regulatory assessment time by 62% and implementation costs by 41% compared to custom implementations [11].

As these frameworks continue to mature, platform providers are developing compliance documentation and audit capabilities that simplify governance for organizations without specialized expertise. These features transform what was previously a barrier to adoption into a manageable operational consideration, further democratizing access across industries and organization sizes while ensuring responsible innovation.
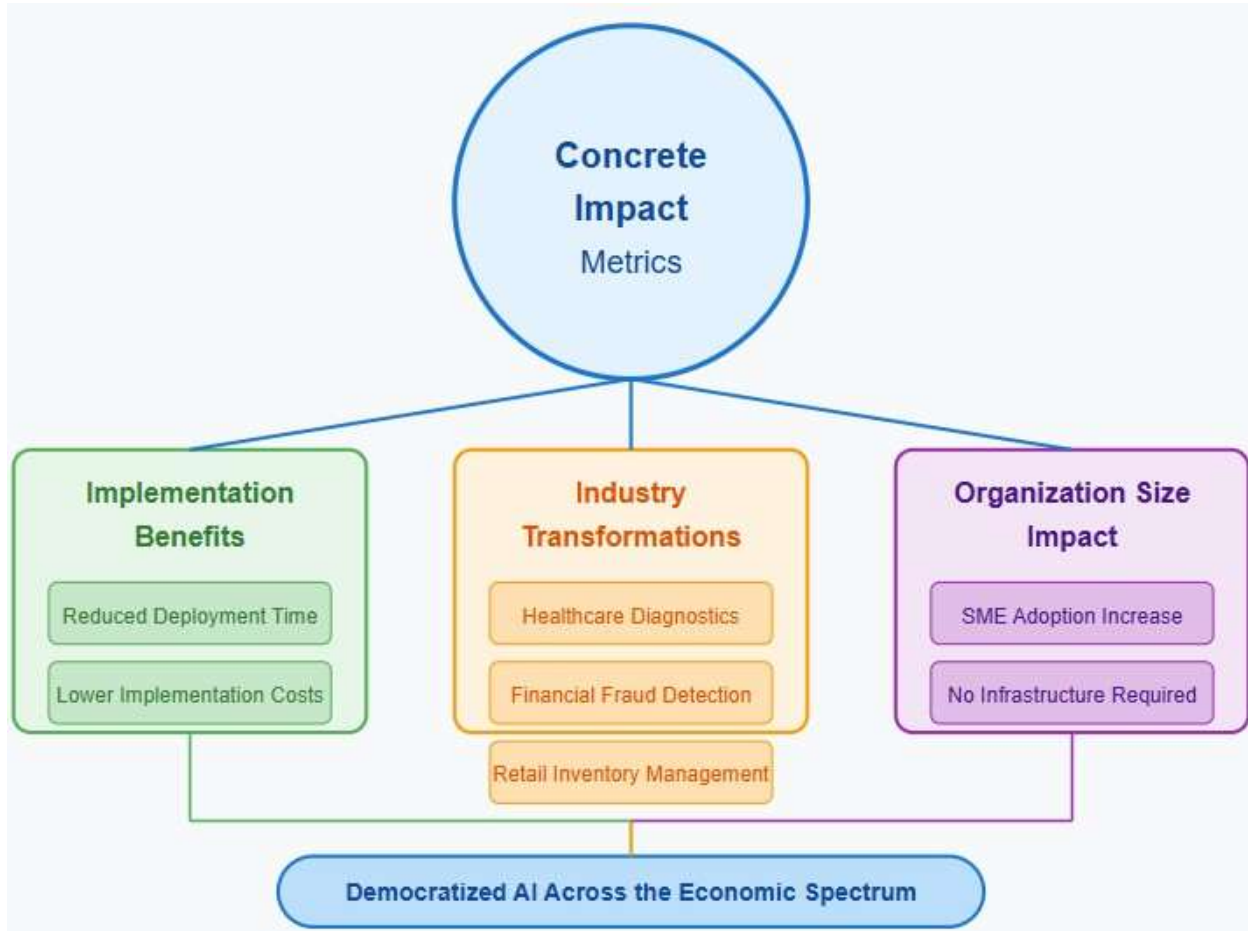
Fig 3: Concrete Impact of AI Democratization Across Industries [9,10]

## 6. Conclusion

The transformation of Generative AI inferencing through hardware acceleration, software optimization, and deployment innovation has fundamentally altered how organizations implement large language models and other generative systems. Modern platforms with standardized APIs and flexible deployment options have democratized access to sophisticated capabilities previously restricted to resource-rich organizations. This democratization has delivered measurable business value across healthcare diagnostics, financial fraud detection, retail operations, and manufacturing optimization. The integration of ethical safeguards and governance frameworks addresses critical considerations around privacy, bias, and regulatory compliance, ensuring responsible implementation. Organizations of all sizes now address specific challenges through accessible platforms rather than custom infrastructure, with documented improvements in implementation speed, cost efficiency, and operational performance. As platforms evolve, simplified interfaces, automated optimization, and consumption-based pricing will accelerate adoption across more organizations and use cases. This shift represents a fundamental change in how Generative AI capabilities are consumed—from specialized technology requiring substantial investment to accessible tools addressing specific business needs. The significance extends beyond technological advancement to economic transformation, creating a more inclusive AI ecosystem that distributes capabilities across the economic spectrum while maintaining appropriate guardrails for responsible use.

**Conflicts of Interest:** The authors declare no conflict of interest.
**Publisher's Note:** All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

## References

[1]   Amna B et al., (2025) AI governance: a systematic literature review, Springer Nature Link, Volume 5, pages 3265–3279, 2025. [Online]. Available: https://link.springer.com/article/10.1007/s43681-024-00653-w

[2]     Andreas S, (2019) The Democratization of Artificial Intelligence: Net Politics in the Era of Learning Algorithms, The Deutsche Nationalbibliothek, 2019. [Online]. Available: https://library.oapen.org/bitstream/id/61810f6f-bf40-404d-b688-61e9ac3405a4/external_content.pdf

[3]     Badahun K et al., (2024) The Evolution of Edge Computing in a Data-Driven World, *International Research Journal on Advanced Engineering and Management (IRJAEM)* 2(10):3238-3250, 2024. [Online]. Available: https://www.researchgate.net/publication/385260179_The_Evolution_of_Edge_Computing_in_a_Data-Driven_World

[4]     Evidently AI, (2025) What is data drift in ML, and how to detect and handle it, 2025. [Online]. Available: https://www.evidentlyai.com/ml-in-production/data-drift

[5]     Grand View Research, (2025) Artificial Intelligence Market Size, Share & Trends Analysis Report By Solution, By Technology (Deep Learning, Machine Learning, NLP, Machine Vision, Generative AI), By Function, By End-Use, By Region, And Segment Forecasts, 2025 - 2030. [Online]. Available: https://www.grandviewresearch.com/industry-analysis/artificial-intelligence-ai-market

[6]     James P, (2025) Privacy-Preserving Techniques in Generative AI Models for Secure Cloud Data Processing, ResearchGate, 2025. [Online]. Available: https://www.researchgate.net/publication/390742972_Privacy-Preserving_Techniques_in_Generative_AI_Models_for_Secure_Cloud_Data_Processing

[7]     Martins A and Kayode S, (2024) Evaluating the Trade-Offs: Cost vs. Performance in Serverless Computing for AI and ML Workload Deployment, ResearchGate, 2024. [Online]. Available: https://www.researchgate.net/publication/386424516_Evaluating_the_Trade-Offs_Cost_vs_Performance_in_Serverless_Computing_for_AI_and_ML_Workload_Deployment

[8]     Mercia L et al., (2025) Exploring Model Compression Techniques for Efficient Inference of Large Language Models, HAL Open Science, 2025. [Online]. Available: https://hal.science/hal-04997150/document#:~:text=Model%20compression%20has%20emerged%20as,%2Drank%20factorization%20%5B109%5D.
Siyuan L et al., (2025) Adaptive lightweight network construction method for Self-Knowledge Distillation, Neurocomputing, Volume 624, 129477, 2025. [Online]. Available: https://www.sciencedirect.com/science/article/abs/pii/S0925231225001493

[9]     Narges L et al., (2024) HPC-AI benchmarks - A comparative overview of high-performance computing hardware and AI benchmarks across domains, *Journal of Artificial Intelligence and Robotics*, Vol. 1, Issue. 1, 2024. [Online]. Available: https://joaiar.org/articles/AIR-1017.pdf

[10]   Sciforce, (2024) MLOps as The Key to Efficient AI Model Deployment and Maximum ROI, Medium, 2024. [Online]. Available: https://medium.com/sciforce/mlops-as-the-key-to-efficient-ai-model-deployment-and-maximum-roi-d9b96cea8412

[11]   Sudhi S & Young M. L, (2024) Challenges with developing and deploying AI models and applications in industrial systems, Volume 4, article number 55, Springer Nature, 2024. [Online]. Available: https://link.springer.com/article/10.1007/s44163-024-00151-2