# Pricing Models for Agentic AI Systems: From Token-Based to Value-Based Approaches

**AI-Generated Academic Thesis Showcase**

Academic Thesis AI (Multi-Agent System)

January 2025

# Table of Contents

# Abstract

**Research Problem and Approach:** The rapid proliferation of agentic AI systems, particularly Large Language Models (LLMs), presents a significant challenge for traditional pricing models, which struggle to account for their dynamic capabilities, variable resource consumption, and often opaque value generation. This thesis addresses the gap in understanding effective monetization strategies for these advanced AI services, moving beyond simplistic token-based approaches. It adopts a comprehensive theoretical analysis, complemented by real-world case studies, to dissect the economic, strategic, and ethical implications of current and emergent AI pricing paradigms.

**Methodology and Findings:** Employing a multi-dimensional theoretical framework, the study systematically compares token-based, usage-based, and value-based pricing models across dimensions such as cost predictability, scalability, market segmentation, and ethical considerations. Key findings reveal that while consumption-based models offer granularity and cost alignment, they often fall short in capturing the full economic value delivered by AI. The research identifies a clear trend towards hybrid models, which combine elements of different paradigms to optimize for diverse user segments and balance provider revenue stability with user cost predictability. Furthermore, it highlights the increasing potential for AI-driven dynamic pricing and the emergent risk of algorithmic collusion, underscoring the need for robust governance.

**Key Contributions:** (1) Development of a multi-dimensional theoretical framework for comparing AI pricing models, integrating economic principles with AI-specific characteristics. (2) Comprehensive analysis of real-world AI pricing implementations by leading providers, detailing their strategic rationales and market implications. (3) Identification and exploration of hybrid pricing typologies, demonstrating their strategic imperative in balancing revenue, adoption, and customer value.

1

**Implications:** This research offers critical insights for AI companies to design sustainable and equitable monetization strategies, for customers to navigate the complexities of AI service costs, and for policymakers to develop appropriate regulatory frameworks. It emphasizes the need for transparency, value articulation, and ethical considerations in AI pricing to foster widespread adoption and ensure fair competition in the evolving digital economy.

**Keywords:** AI pricing, LLM pricing, Agentic AI, Token-based pricing, Usage-based pricing, Value-based pricing, Hybrid pricing models, Algorithmic collusion, Dynamic pricing, Market segmentation, Cloud services, API monetization, Pricing strategy, Economic models, AI governance

## Introduction

The arrival of artificial intelligence (AI) signals a huge shift. It's fundamentally reshaping industries, economies, and even society–and doing so at an unprecedented pace (Korinek, 2025)(Geetha et al., 2024). From automating complex tasks to generating new insights, AI's capabilities are growing fast. Large language models (LLMs) and, more recently, agentic AI systems are emerging as real game-changers. These advanced systems aren't just tools. Instead, they're becoming more autonomous, able to reason, plan, and execute complex tasks across many fields (Barbere et al., 2024). As advanced AI becomes more integrated into businesses and consumer apps, how we manage its economics–especially pricing–becomes a key factor. This impacts widespread adoption, fair access, and sustainable development (Lorente, 2025)(Brynjolfsson et al., 2023). Figuring out how to price AI, particularly for dynamic, autonomous agentic systems, is quickly becoming a major problem for tech providers, platform operators, and consumers alike.

Traditional software pricing–think subscriptions, per-user licenses, or fixed tiers–doesn't quite fit the complex value and changing resource use of modern AI applications (De, 2017)(Seufert, 2014). This complexity is made worse by AI's often unclear operations, unpredictable results, and continuously evolving models. For agentic AI, these challenges are even bigger due to their inherent autonomy and goal-oriented behaviors. Unlike standard API calls that do one specific thing, agentic AI can orchestrate multi-step processes, interact with various external tools, and adapt strategies in real-time to achieve objectives. This dynamic, often unpredictable operational footprint makes pricing them uniquely challenging.

## Literature Review

The rapid proliferation of Artificial Intelligence (AI) and Large Language Models (LLMs) has inaugurated a new era of digital services, presenting both unprecedented opportunities and complex challenges for their economic valuation and monetization. As these

3

technologies become increasingly embedded in diverse industries, from healthcare and finance to automotive and retail, the strategies for pricing and cost optimization have evolved significantly, moving beyond traditional software and cloud service models (Shiva Kumar Bhuram, 2025)(Subham, 2025). This literature review synthesizes existing research on AI and LLM pricing, drawing on foundational economic theories, digital goods pricing models, and specific strategies adopted by leading AI service providers. It aims to provide a comprehensive understanding of token-based pricing, usage-based pricing, and value-based pricing, culminating in a comparative analysis and an exploration of emerging trends and their implications. The discussion will highlight how these pricing paradigms attempt to reconcile the unique characteristics of AI services–such as their high initial development costs, near-zero marginal replication costs, dynamic capabilities, and often opaque operational expenses–with the imperative for sustainable business models and equitable access.

The economic landscape of AI services is distinguished by several key factors that complicate conventional pricing approaches. Unlike tangible goods, AI models are non-rivalrous and often non-excludable once developed, creating challenges akin to public goods or information goods (Fishburn & Odlyzko, 1999). The cost structure is heavily front-loaded, with substantial investments in research, data acquisition, model training, and infrastructure, yet the marginal cost of serving an additional user or generating another output can be negligible (Cho & Bahn, 2020). This cost asymmetry necessitates pricing models that can recoup initial investments while remaining competitive and scalable. Furthermore, the value derived from AI is often highly contextual and user-dependent, ranging from direct efficiency gains and cost reductions to enhanced decision-making capabilities and improved customer experiences (Lorente, 2025)(Maguire, 2021). Understanding these intricate dynamics is crucial for developing robust pricing strategies that align with both provider sustainability and user utility.

The evolution of pricing in the digital realm provides a critical backdrop for understanding current AI monetization strategies. Early digital products often adopted subscription

or license-based models, mirroring traditional software (Fishburn & Odlyzko, 1999). However, the advent of cloud computing fundamentally shifted this paradigm towards more granular, usage-based approaches, where customers pay for what they consume, be it compute cycles, storage, or network bandwidth (Livingstone, 2013)(Cho & Bahn, 2020). This shift was driven by the desire for greater flexibility, scalability, and cost alignment for users, allowing them to scale resources up or down dynamically. AI services, particularly those delivered via Application Programming Interfaces (APIs), inherit many of these characteristics but introduce new complexities related to the nature of AI processing, such as the concept of "tokens" in LLMs or the specific computational demands of various AI tasks (Barbere et al., 2024)(De, 2017). This review will delve into how these different layers of pricing theory and practice are being adapted and innovated to address the unique economic footprint of AI.

The literature review is structured to first establish the foundational economic principles relevant to pricing digital and cloud services, laying the groundwork for understanding the specificities of AI. Following this, it will critically examine token-based pricing models, which have emerged as a dominant method for monetizing LLMs, discussing their mechanics, advantages, and inherent challenges. Subsequently, the review will explore broader usage-based pricing models, encompassing not only cloud infrastructure but also API-driven AI services, analyzing their implementation, benefits, and complexities. A dedicated section will then address value-based pricing theory, exploring its application in the context of AI, where the perceived utility and impact on business outcomes often outweigh direct computational costs. Finally, a comparative analysis will synthesize these different approaches, highlighting hybrid models, the role of dynamic pricing, and the broader competitive and regulatory landscape, before outlining future research directions.

*Foundational Theories of Pricing in Digital and Cloud Services*

The pricing of AI and LLM services is deeply rooted in established economic theories, albeit with significant adaptations required for the unique characteristics of digital goods and

services. Traditional pricing theory often begins with cost-plus pricing, where a markup is added to production costs, or competitive pricing, where prices are set in relation to market rivals (Lo, 2018). However, digital goods, including AI models, fundamentally challenge these traditional frameworks due to their distinct cost structures and consumption patterns. The marginal cost of producing an additional copy or instance of a digital good approaches zero after the initial development (Fishburn & Odlyzko, 1999). This "first-copy problem" means that recouping high fixed costs requires different strategies than those for physical goods with significant variable costs.

One key theoretical lens is the economics of information goods, which highlights characteristics such as non-rivalry (one person's consumption does not diminish another's) and often low excludability (Fishburn & Odlyzko, 1999). For AI services, this translates to the ability to serve numerous users with the same underlying model, making scale a critical factor in profitability. Fishburn and Odlyzko (1999) extensively discuss competitive pricing for information goods, noting the prevalence of subscription pricing and versioning strategies to differentiate offerings and cater to diverse customer segments (Fishburn & Odlyzko, 1999). Versioning, which involves offering different qualities or feature sets at various price points, is particularly relevant for AI, where models can be offered in different sizes, with varying performance levels, or with specialized capabilities (Diaw & Pouyet, 2004). This allows providers to capture surplus from different market segments by offering a range of price-performance trade-offs (Sonderegger, 2011).

The advent of cloud computing further revolutionized pricing models by shifting from traditional software licensing to a utility-based, pay-per-use paradigm (Ladas et al., 2019). Livingstone (2013) coined the term "volatility as a service," emphasizing the dynamic nature of cloud pricing and resource allocation (Livingstone, 2013). This model allows users to scale their infrastructure up or down based on demand, paying only for the resources consumed (Cho & Bahn, 2020). This flexibility is a significant advantage for businesses with fluctuating workloads, as it eliminates the need for large upfront capital expenditures on hardware. Cloud

providers like AWS, Azure, and Google Cloud have developed sophisticated usage-based pricing structures that encompass compute, storage, networking, and various specialized services (Satapathi, 2025). Cho and Bahn (2020) proposed a cost estimation model for cloud services, highlighting the complexity of accurately predicting expenses in a pay-as-you-go environment (Cho & Bahn, 2020). This complexity stems from the numerous variables involved, including data transfer costs, different instance types, and varying regional prices, making cost optimization a significant challenge for users (Anonymous, 2025).

Market segmentation plays a crucial role in digital and cloud service pricing. Sonderegger (2011) examined market segmentation with nonlinear pricing, where the price per unit changes based on the quantity consumed (Sonderegger, 2011). This strategy is evident in tiered pricing models, volume discounts, and different subscription levels that offer varying allowances of usage. Nonlinear pricing allows providers to extract more value from high-volume users while still attracting low-volume users, effectively segmenting the market based on demand elasticity (Ye & Zhang, 2017). In the context of AI, this can manifest as different pricing for various tiers of API calls, token usage, or access to advanced features. For instance, free tiers, as discussed by Seufert (2014) in the context of freemium products, serve as a powerful customer acquisition tool, allowing users to experience the service before committing to paid usage (Seufert, 2014).

Transaction Cost Economics (TCE), as articulated by Williamson (2010), also offers insights into AI service pricing (Williamson, 2010). TCE focuses on the costs associated with making economic exchanges, including search and information costs, bargaining costs, and enforcement costs. In the context of AI, leveraging external AI services via APIs rather than developing capabilities in-house can reduce transaction costs related to development, maintenance, and scaling. The pricing model, therefore, influences the "make or buy" decision for organizations. A transparent and predictable pricing model can reduce perceived transaction risks, encouraging adoption. Conversely, opaque or highly variable pricing

can increase perceived transaction costs, disincentivizing reliance on external AI providers (Mirghaderi et al., 2023).

The competitive landscape also significantly influences pricing strategies. Cody (2000) discussed competitive infrastructure as an enabler of market-based pricing, where robust competition among providers drives down prices and encourages innovation (Cody, 2000). The "cloud price wars" described by Livingstone (2013) exemplify this phenomenon, where major cloud providers aggressively cut prices to gain market share (Livingstone, 2013). This competitive pressure is now extending to the AI service market, as numerous providers vie for dominance, leading to downward pressure on prices and the proliferation of more sophisticated pricing models (Liu, 2024). However, concerns about market power and potential for excessive pricing by dominant players remain, particularly given the network effects and data moats often associated with large AI platforms (Ayata, 2020)(Tucker, 2019). The regulatory environment and antitrust considerations, as highlighted by Tucker (2019) and Lo (2018), play a crucial role in shaping these competitive dynamics and ensuring fair pricing practices (Tucker, 2019)(Lo, 2018).

*Token-Based Pricing Models in Large Language Models (LLMs)*

The emergence of Large Language Models (LLMs) has introduced a novel pricing paradigm centered around the concept of "tokens." Token-based pricing is a granular, usage-based model specifically tailored to the discrete units of text processing performed by LLMs. A "token" is a sequence of characters, often a word or sub-word unit, into which an LLM breaks down input text (prompts) and generates output text (completions) (Barbere et al., 2024). This approach directly ties the cost of an LLM interaction to the computational load it incurs, as processing more tokens generally requires more computational resources. The widespread adoption of token-based pricing by leading LLM providers like OpenAI and Anthropic underscores its perceived efficiency and fairness in a domain where computational demands vary wildly across different queries.

The mechanics of tokenization are complex and vary slightly between models, impacting how users are charged. Generally, a token can be a word, a punctuation mark, or even parts of a word, especially for less common terms (Barbere et al., 2024). For example, the word "tokenization" might be broken into "token," "iza," and "tion," each counting as a separate token. This means that the number of tokens is not always directly equivalent to the number of words, making cost estimation challenging for end-users. Barbere et al. (2024) discuss dynamic token hierarchies, suggesting an evolution in how models process and potentially price these units, moving towards more nuanced representations that could affect future pricing structures (Barbere et al., 2024). Rudnytskyi (2022) provides an R wrapper for the OpenAI API, implicitly demonstrating the programmatic interaction required to manage token usage (Rudnytskyi, 2022).

One of the primary advantages of token-based pricing is its high granularity. Users are charged precisely for the input they provide and the output they receive, which can be seen as fair, especially for diverse applications. For tasks requiring short, precise answers, token counts remain low, resulting in minimal costs. Conversely, for extensive content generation or summarization tasks, where both input prompts and output completions can be lengthy, the costs scale proportionally (Barbere et al., 2024). This direct correlation with usage makes it a flexible model for developers integrating LLMs into applications with varying demands. Moreover, token-based pricing aligns with the underlying computational reality of LLMs, where the processing of each token consumes a measurable amount of computational power, including CPU, GPU, and memory resources.

However, token-based pricing is not without its challenges and criticisms. A significant hurdle for users is the lack of transparency and predictability in cost estimation. Since tokenization rules are internal to each model and can vary, it is difficult for users to accurately predict how many tokens a given piece of text will consume without first running it through the model's tokenizer (Barbere et al., 2024). This unpredictability can lead to "bill shock" for applications that generate unexpected volumes of text or process large inputs. Furthermore,

LLMs often price input tokens differently from output tokens, typically with output tokens being more expensive due to the generation process being more computationally intensive (Barbere et al., 2024). This differential pricing adds another layer of complexity for budgeting and cost control.

The concept of the "context window" is intimately linked to token-based pricing. LLMs have a finite context window, which is the maximum number of tokens they can process at once, including both the input prompt and the generated output (Barbere et al., 2024). As models evolve to support larger context windows, the potential for higher token usage and thus higher costs increases. While larger context windows enable more sophisticated applications, such as processing entire documents or maintaining long-running conversations, they also demand more computational resources per inference, potentially leading to higher per-token costs or increased overall expenditure for users who leverage these capabilities extensively. Barbere et al. (2024) explore how dynamic token hierarchies could enhance LLMs, implying that future models might manage and price tokens in more sophisticated ways, potentially optimizing for both performance and cost (Barbere et al., 2024).

Users and developers have adopted various strategies to optimize token usage and manage costs. Prompt engineering, for instance, involves crafting concise yet effective prompts to minimize input token count while still eliciting desired responses (Barbere et al., 2024). Techniques such as summarization of prior conversations or selective retrieval of information can also reduce the amount of data fed into the LLM, thereby lowering token consumption. For applications requiring extensive context, developers might employ retrieval-augmented generation (RAG) architectures, where relevant information is retrieved from a database and dynamically inserted into the prompt, rather than feeding the entire database to the LLM, thus managing token limits and costs more effectively.

Ethical implications also arise in token-based pricing. The opacity of tokenization and the potential for providers to adjust token counting mechanisms could lead to concerns about fairness and transparency (Mirghaderi et al., 2023). If users cannot independently

verify token counts, trust in the pricing model may erode. Moreover, the incentive structure created by token-based pricing might encourage providers to optimize for token generation rather than conciseness or efficiency in certain scenarios, although this is generally mitigated by competitive pressures and user demand for economical solutions. The balance between providing powerful, flexible models and ensuring transparent, predictable costs remains a critical challenge for LLM providers.

In summary, token-based pricing represents a direct and granular approach to monetizing LLMs, aligning costs with the fundamental units of processing in these models. Its advantages lie in its flexibility and direct correlation to computational load, making it suitable for a wide range of applications. However, the inherent complexities of tokenization, cost predictability, and the interplay with context window sizes present significant challenges for users seeking to manage and optimize their AI expenditures. As LLMs continue to evolve, so too will the sophistication of token-based pricing, potentially incorporating dynamic hierarchies and more transparent mechanisms to address current limitations (Barbere et al., 2024).

*Usage-Based Pricing Models in AI and Cloud Services*

Usage-based pricing models extend beyond the specific token counting of LLMs to encompass a broader spectrum of AI and cloud services, where customers are charged based on their consumption of various resources or features. This paradigm has been a cornerstone of cloud computing since its inception and has profoundly influenced how AI services are delivered and monetized (Livingstone, 2013). Unlike subscription or license-based models that charge a fixed fee regardless of usage, usage-based pricing offers flexibility, scalability, and a direct alignment of costs with actual consumption, making it particularly attractive for dynamic workloads and varying user demands (Ladas et al., 2019).

The foundational principles of usage-based pricing are deeply embedded in the offerings of major cloud providers like Amazon Web Services (AWS), Microsoft Azure, and Google

Cloud Platform (GCP). These platforms charge for a myriad of resources, including compute instances (e.g., CPU hours, GPU hours), storage (e.g., gigabytes per month), data transfer (e.g., egress bandwidth), and specific managed services (Satapathi, 2025). This granular billing allows businesses to scale their infrastructure up or down on demand, paying only for what they use, which significantly reduces upfront capital expenditures and operational overhead. Cho and Bahn (2020) analyzed cost estimation models for cloud services, underscoring the complexity arising from the multitude of pricing dimensions and the dynamic nature of resource consumption (Cho & Bahn, 2020). Despite the complexity, the flexibility offered by these models has been a key driver of cloud adoption (Livingstone, 2013).

In the realm of AI, usage-based pricing manifests in several forms. Beyond the token-based models for LLMs, other AI services are typically priced per API call, per unit of data processed, or per specific feature invoked. For example, a sentiment analysis API might charge per text analyzed, an image recognition service per image processed, or a speech-to-text service per minute of audio transcribed (Leechewyuwasorn & Wangpratham, 2024)(Satapathi, 2025). This "API monetization" strategy, as discussed by De (2017), is a common approach for exposing AI capabilities as modular services (De, 2017). Rodeghero et al. (2017) highlighted the importance of API usage in descriptions of source code functionality, indicating that the API is often the primary interface through which developers interact with and are charged for AI services (Rodeghero et al., 2017).

The advantages of usage-based pricing are multifaceted. Firstly, it provides unparalleled flexibility and scalability. Users can start small with minimal investment and scale up as their needs grow, making it accessible for startups and large enterprises alike. Secondly, it aligns costs directly with actual consumption, which can lead to greater cost efficiency if usage is managed effectively. This pay-as-you-go model minimizes waste, as resources are only provisioned when needed. Thirdly, it lowers the barrier to entry for adopting advanced AI technologies, as users do not need to commit to expensive licenses or infrastructure upfront (Seufert, 2014). This democratizes access to powerful AI tools, enabling broader innovation.

However, usage-based pricing also presents significant challenges, particularly concerning cost management and predictability. The sheer number of variables and the dynamic nature of consumption can make it extremely difficult for organizations to forecast and control their AI and cloud spending, often leading to unexpected "bill shock" (Anonymous, 2025). Users must invest in sophisticated cost monitoring and optimization tools to track usage across various services and implement strategies to minimize unnecessary expenses. Anonymous (2025) highlights AI-driven strategies for cloud cost optimization, suggesting that AI itself can be leveraged to manage and predict costs associated with its own consumption (Anonymous, 2025). This creates a fascinating feedback loop where AI helps optimize the cost of AI.

Another challenge is vendor lock-in. While usage-based models offer flexibility in scaling, switching providers can be costly due due to data migration efforts, re-architecting applications, and retraining personnel (Vogt & Bizer, 2013). This can limit competition and potentially allow dominant providers to exert market power over time, leading to concerns about excessive pricing (Ayata, 2020). Kaaniche and Laurent (2018) addressed blockchain-based data usage auditing, which could potentially enhance transparency and accountability in usage tracking, mitigating some of these lock-in concerns by providing verifiable records of consumption (Kaaniche & Laurent, 2018).

Tiered pricing models are a common variant of usage-based pricing, offering different price points based on volume or specific feature sets. These often include free tiers for initial exploration, standard pay-as-you-go rates, and discounted rates for committed usage or higher volumes (Seufert, 2014)(Satapathi, 2025). Satapathi (2025) outlines pricing tiers for Azure AI Language Service, illustrating how different levels of service or specific functionalities within an AI suite are priced (Satapathi, 2025). This allows providers to cater to diverse customer segments, from individual developers to large enterprises, optimizing revenue capture across the demand curve (Sonderegger, 2011).

The integration of AI into edge computing environments further complicates usage-based pricing. Shiva Kumar Bhuram (2025) discusses Edge-Cloud AI for dynamic pricing in automotive aftermarkets, indicating that AI processing can occur closer to the data source, reducing latency and bandwidth costs (Shiva Kumar Bhuram, 2025). Pricing models for such hybrid deployments must account for the distributed nature of computation, potentially involving charges for local processing, data synchronization, and cloud-based orchestration. This highlights the need for flexible and adaptive pricing mechanisms that can span heterogeneous computing environments.

Furthermore, "Green AI" considerations are increasingly influencing usage-based pricing. Kshirsagar et al. (2021) propose Green Artificial Intelligence Powered Cost Pricing Models for Congestion Control (GREE-COCO), which integrates environmental sustainability objectives alongside cost optimization (Kshirsagar et et al., 2021). This approach suggests that future usage-based models might incorporate factors related to energy consumption and carbon footprint, incentivizing more energy-efficient AI deployments and usage patterns (Gbadamosi et al., 2018). Such models would not only reflect computational costs but also the ecological impact of AI processing, moving towards a more holistic pricing framework.

In conclusion, usage-based pricing models, inherited from the cloud computing paradigm, are fundamental to the monetization of AI services. They offer significant advantages in flexibility, scalability, and cost alignment with consumption, facilitating broader adoption of AI. However, they also pose substantial challenges related to cost predictability, management complexity, and potential vendor lock-in. The ongoing evolution includes sophisticated tiered structures, adaptations for edge-cloud deployments, and the nascent integration of environmental sustainability metrics, all aimed at refining how AI services are economically valued and delivered (Kshirsagar et al., 2021)(Shiva Kumar Bhuram, 2025)(Anonymous, 2025).

*Value-Based Pricing Theory in AI Services*

Beyond the quantitative metrics of tokens or raw usage, value-based pricing theory offers a distinct paradigm for monetizing AI services, focusing on the perceived or actual economic value delivered to the customer (Maguire, 2021). This approach represents a shift from cost-centric or competitor-centric pricing to a customer-centric model, where the price is determined by the benefits and outcomes that the AI solution enables for the user (Lorente, 2025). While more complex to implement, value-based pricing holds significant promise for AI, given its transformative potential across various business functions.

The rationale for value-based pricing in AI stems from the often profound impact these technologies can have on an organization's bottom line or strategic objectives. AI solutions can drive revenue growth through personalized product bundling (Chiruvelli, 2025) and optimized sales forecasting (Subham, 2025), reduce operational costs through improved efficiency and automation (Abbas, 2025), enhance decision-making with predictive analytics (Niharika et al., 2024), or create entirely new business models (Lorente, 2025). Maguire (2021) emphasizes "value selling," where the focus is on articulating and demonstrating the quantifiable benefits a product or service provides to the customer, rather than merely listing features or technical specifications (Maguire, 2021). Lorente (2025) further explores value creation and value capture in AI through a Triple Helix Model, highlighting the interplay between academia, industry, and government in realizing AI's economic potential (Lorente, 2025).

Measuring and quantifying the value of AI, however, presents a significant challenge. Unlike a tangible product with a clear purchase price and resale value, the value of an AI service can be diffuse, indirect, and highly context-dependent. It often requires a deep understanding of the customer's business processes, strategic goals, and key performance indicators (KPIs) (Maguire, 2021). Tangible benefits might include a measurable increase in sales conversion rates (Subham, 2025), a reduction in customer churn, or significant time savings from automating routine tasks (Abbas, 2025). Intangible benefits, while harder

to quantify, are equally important and can include improved brand reputation, enhanced customer experience (Fang & Zhou, 2025), better employee satisfaction, or increased agility in responding to market changes (Siddannavar et al., 2025). Fang and Zhou (2025) delve into understanding the impacts of human-like competencies on user engagement, which indirectly contributes to the perceived value of AI services (Fang & Zhou, 2025).

One of the primary difficulties in applying value-based pricing to AI lies in isolating the specific contribution of the AI solution from other factors influencing business outcomes. For instance, if an AI-powered recommendation system leads to a 10% increase in e-commerce sales, how much of that increase is attributable solely to the AI, and how much to marketing efforts, product quality improvements, or seasonal demand? This attribution problem requires robust methodologies for baseline measurement, A/B testing, and careful analytical frameworks to demonstrate the direct impact of the AI (Maguire, 2021). This complexity often necessitates a consultative sales approach, where AI providers work closely with clients to define success metrics and establish clear value propositions (Oleksii, 2025).

Ethical considerations also play a role in value-based pricing for AI. If AI models are used in critical domains, such as healthcare or finance, their value might be tied to outcomes that have significant societal implications. Mirghaderi et al. (2023) discuss ethics and transparency issues in digital platforms, which are highly relevant when pricing is based on the perceived value of potentially biased or opaque AI systems (Mirghaderi et al., 2023). Ensuring that value-based pricing does not inadvertently lead to discriminatory outcomes or exploit vulnerable populations is paramount. The concept of "fairness" in pricing, especially in markets with potential for market power, remains a contentious issue (Ayata, 2020).

Strategies for implementing value-based pricing in AI services often involve outcome-based models, where payment is directly tied to the achievement of specific business results. For example, an AI provider might charge a percentage of the revenue increase generated by their recommendation engine or a fixed fee per successful fraud detection. This approach aligns the interests of the provider and the customer, as the provider's compensation is

directly linked to the value they create (Ladas et al., 2019). However, outcome-based pricing requires robust contracts, clear definitions of success, and mechanisms for verifying outcomes, which can add complexity to the commercial relationship.

Another strategy involves offering tiered value packages, where different service levels are priced based on the depth of insights, level of automation, or guaranteed performance metrics. For instance, a basic AI analytics package might offer descriptive insights, while a premium package provides predictive capabilities, prescriptive recommendations, and higher service level agreements (SLAs). This allows customers to choose a tier that best matches their perceived value and budget (Sonderegger, 2011). Chaerul (2025) emphasizes analyzing user expectation and experience to formulate data-driven pricing strategies, which is crucial for aligning tiered offerings with customer value perception (Chaerul, 2025).

The ability of AI itself to assist in value assessment and pricing optimization is an emerging area. Korinek (2025) explores the use of AI agents for economic research, suggesting that AI could be employed to analyze market data, predict demand elasticity, and even recommend optimal pricing strategies based on value propositions (Korinek, 2025). Oleksii (2025) discusses algorithmizing B2B sales, where AI can help create sales frameworks that incorporate value-based arguments and optimize pricing negotiations (Oleksii, 2025). This represents a fascinating feedback loop where AI tools are used to refine the monetization strategies for AI services themselves.

Moreover, the psychological factors affecting customer lifetime value, as analyzed by Siddannavar et al. (2025), are intrinsically linked to value perception (Siddannavar et al., 2025). If customers perceive high value from an AI service, their loyalty and willingness to pay will increase, contributing to long-term profitability. Yin and Qiu (2021) explored AI technology and online purchase intention, highlighting how the perceived utility and impact of AI can directly influence consumer behavior and willingness to adopt (Yin & Qiu, 2021). This suggests that demonstrating and communicating value effectively is not just a pricing strategy but a fundamental aspect of market adoption.

In conclusion, value-based pricing, while conceptually appealing for AI services, requires overcoming significant hurdles in quantifying, attributing, and communicating the value generated. It necessitates a deep understanding of customer needs, robust measurement frameworks, and often a consultative approach. Despite these complexities, its potential to align provider incentives with customer outcomes, especially in transformative AI applications, makes it a critical component of a comprehensive AI monetization strategy (Maguire, 2021)(Lorente, 2025). As AI becomes more sophisticated and its impacts more profound, the ability to articulate and price based on delivered value will become increasingly important for sustainable growth and market differentiation.

*Comparative Analysis and Emerging Trends*

The landscape of AI and LLM pricing is characterized by a dynamic interplay between token-based, usage-based, and value-based models. While each approach has distinct merits and challenges, the market is increasingly witnessing the emergence of hybrid models that combine elements from multiple paradigms to create more flexible, comprehensive, and competitive offerings. This comparative analysis synthesizes these approaches, identifies critical emerging trends such as dynamic pricing and the influence of competition and regulation, and outlines future directions for AI monetization.

Token-based pricing, primarily for LLMs, offers granular control and a direct link to computational processing, making it transparent in terms of resource consumption (Barbere et al., 2024). Its strength lies in its ability to handle highly variable workloads, where the cost scales precisely with the amount of text processed. However, its main drawbacks are the unpredictability of token counts for users and the potential for opaque tokenization rules. Usage-based pricing, a broader category encompassing API calls, compute time, and data processed, provides similar benefits of flexibility and scalability, derived from its cloud computing heritage (Satapathi, 2025)(Livingstone, 2013). It lowers entry barriers but can lead to complex cost management and "bill shock" if not carefully monitored (Anonymous,

2025). Both token-based and general usage-based models are fundamentally cost-oriented or resource-oriented, tying price directly to the provider's operational expenses or the raw consumption of resources.

In contrast, value-based pricing shifts the focus from inputs to outcomes, aiming to capture the economic benefits an AI solution delivers to the customer (Maguire, 2021). This approach is theoretically ideal for highly impactful AI applications but is challenging to implement due to difficulties in quantifying and attributing value (Lorente, 2025). Its advantage lies in aligning the provider's incentives with the customer's success, potentially leading to higher revenue capture for transformative solutions. However, it requires extensive collaboration, robust measurement frameworks, and can be perceived as less transparent than direct usage-based billing (Mirghaderi et al., 2023).

The current trend leans heavily towards **hybrid models** that strategically combine these approaches. For instance, an LLM service might offer a base subscription (value-based) for a certain number of tokens (token-based), with additional tokens charged on a pay-per-use basis (Barbere et al., 2024)(Seufert, 2014). Furthermore, premium features or higher performance tiers might be priced based on the enhanced value they provide, layered on top of a usage-based foundation (Sonderegger, 2011). This allows providers to offer predictable costs for basic usage while capturing additional revenue from high-value or high-volume consumption. These hybrid models attempt to balance the transparency and scalability of usage-based pricing with the revenue optimization potential of value-based approaches.

**Table 1: Comparative Analysis of Digital and Cloud Service Pricing Models**

| Model Type | Primary Metric | Key Advantage | Main Disadvantage | AI Relevance |
|---|---|---|---|---|
| **Subscription** | Fixed Fee | Predictable cost | Under/Over-utilization | Basic LLM access, bundled AI features |
| **Usage-Based** | Consumption | Flexible, scalable | Cost unpredictability | API calls, compute hours, data processing |

19

| Model Type | Primary Metric | Key Advantage | Main Disadvantage | AI Relevance |
|---|---|---|---|---|
| **Token-Based** | Tokens consumed | Granular, resource-aligned | Opaque tokenization | LLM input/output, generative AI |
| **Value-Based** | Economic outcome | Value capture, ROI | Value quantification | Specialized AI agents, business solutions |
| **Freemium** | Basic access (free) | Customer acquisition | Conversion challenges | Entry-level AI tools, trial periods |
| **Tiered** | Volume/Features | Market segmentation | Complexity, analysis paralysis | Different LLM models, feature sets |

*Note: This table summarizes the core characteristics and relevance of various pricing models in the context of digital and cloud services, including AI. Hybrid models often combine elements from these categories.*

**Dynamic pricing** is another critical emerging trend, powered by AI itself. This involves adjusting prices in real-time based on factors such as demand, supply, time of day, user segment, and competitive offerings (Niharika et al., 2024)(Vashishtha et al., 2022). Predictive analytics, often AI-driven, can forecast demand and optimize pricing for various services, from sales and revenue operations (Subham, 2025) to inventory management (Vashishtha et al., 2022)(Abbas, 2025). Nagaraju et al. (2023) explored predictive modeling for various applications, underscoring the capability of AI to inform dynamic pricing (Nagaraju et al., 2023). The airline industry, as noted by Divakaruni and Navarro (2024), has long utilized dynamic pricing for technology adoption, a model that AI services can increasingly emulate (Divakaruni & Navarro, 2024). Ma et al. (2015) discussed the re-engineering of procurement through AI, which can include dynamic pricing for inputs (Ma et al., 2015). This real-time adjustment maximizes revenue for providers while potentially offering competitive rates to users during off-peak times. However, dynamic pricing introduces further complexity and

can raise concerns about fairness and transparency if not implemented carefully (Mirghaderi et al., 2023)(Obiajulu et al., 2025). Huang et al. (2024) explored multi-participant double auctions for resource allocation, which represent a sophisticated form of dynamic pricing in competitive environments (Huang et al., 2024).

**Competition and regulation** significantly shape the pricing landscape. As more players enter the AI market, competitive pressures drive innovation in pricing models and often lead to price reductions (Cody, 2000)(Liu, 2024). However, the emergence of dominant platforms with significant market power raises concerns about anti-competitive practices, such as excessive pricing or tying arrangements (Ayata, 2020)(Tucker, 2019). Lo (2018) highlighted the intersection of marketer's pricing strategy and competition law, emphasizing the need for ethical and legal compliance (Lo, 2018). Pricing transparency and innovative pricing models are seen as drivers for competitive markets (Obiajulu et al., 2025). AI governance and responsible AI frameworks, as discussed by Ganguly (2025), are becoming crucial for ensuring fair and ethical pricing practices, especially given the potential for AI to influence pricing decisions (Ganguly, 2025).

The **price elasticity of demand** for AI software is a critical factor for providers. Brynjolfsson et al. (2023) provided evidence on this elasticity, showing how demand for AI software responds to price changes (Brynjolfsson et al., 2023). Understanding this elasticity is crucial for setting optimal prices and predicting market adoption. As AI capabilities become more commoditized, providers will need to continually innovate their pricing strategies to reflect evolving value propositions and competitive dynamics.

Future directions for AI pricing include the development of even more sophisticated, AI-driven pricing engines that can optimize revenue, manage costs, and adapt to market conditions in real-time (Korinek, 2025)(Subham, 2025). This could involve leveraging machine learning to predict customer lifetime value (Siddannavar et al., 2025) and tailor pricing offers accordingly. The integration of "Green AI" principles into pricing models, as proposed by Kshirsagar et al. (2021), will likely gain traction, incentivizing environmentally responsible AI

consumption (Kshirsagar et al., 2021). Furthermore, as AI agents become more autonomous, their own internal cost-benefit analysis and resource allocation decisions could influence the pricing of their services, creating a complex ecosystem of AI-driven economic interactions (Korinek, 2025). The increasing adoption of AI on a global scale, as discussed by Geetha et al. (2024), further emphasizes the need for robust and adaptable pricing strategies that can navigate diverse economic and regulatory environments (Geetha et al., 2024).

In conclusion, AI and LLM pricing is rapidly evolving, moving beyond simple usage metrics to embrace complex hybrid models, dynamic pricing, and value-based approaches. The ongoing challenge for providers is to balance the need for revenue generation and cost recovery with the imperative for transparency, fairness, and customer value. As AI continues to mature and integrate deeper into the economy, the sophistication of its monetization strategies will be a key determinant of its widespread adoption and sustainable impact.

# Methodology

The development of a robust methodology is paramount for systematically analyzing the nascent and rapidly evolving field of artificial intelligence (AI) pricing models. This section delineates the research design, theoretical framework employed for comparison, criteria for case study selection, and the subsequent analytical approach. Given the exploratory nature of AI pricing, which often transcends traditional paradigms, a mixed-methods approach combining theoretical analysis with in-depth case studies is adopted (Sibuea & Arfianti, 2021). This approach facilitates the examination of complex phenomena within their real-world contexts, providing rich, nuanced insights that are critical for understanding the strategic implications of different pricing models (Aberdeen, 2013). The objective is to move beyond mere descriptive accounts to generate actionable insights and contribute to a more comprehensive theoretical understanding of value creation and capture in the AI domain (Lorente, 2025).

*Theoretical Framework for AI Pricing Model Comparison*

To systematically compare diverse AI pricing models, a multi-dimensional theoretical framework is developed, integrating established economic principles with the unique characteristics of AI services. This framework serves as an analytical lens, enabling a structured assessment of how different pricing strategies address value proposition, cost structure, competitive dynamics, user behavior, and ethical considerations. The foundational premise is that effective AI pricing must balance the inherent complexities of AI technology with market demands and strategic business objectives (De, 2017)(Niharika et al., 2024).

**Foundational Economic Principles of Pricing**  Traditional pricing theories provide a critical starting point for understanding AI pricing, even as AI introduces novel complexities. Cost-plus pricing, for instance, involves adding a markup to the direct and indirect costs of production (Nagle & Müller, 2017). While seemingly straightforward, this approach becomes intricate with AI due to the difficulty in attributing specific costs to individual service units, especially with shared infrastructure and continuous learning models (Kshirsagar et al., 2021)(Cho & Bahn, 2020). Value-based pricing, conversely, anchors prices to the perceived value delivered to the customer (Maguire, 2021). This is particularly relevant for AI, where the value can be highly subjective, context-dependent, and may evolve over time as the AI system improves or integrates deeper into customer workflows (Lorente, 2025)(Chaerul, 2025). Competitive pricing, on the other hand, involves setting prices based on competitors' prices (Cody, 2000)(Lo, 2018). In the burgeoning AI market, this often leads to "cloud price wars" (Livingstone, 2013), where providers aggressively lower prices to gain market share, a strategy that may not be sustainable or reflective of the true value and cost of advanced AI capabilities.

Beyond these fundamental approaches, several advanced economic concepts are pertinent. Market segmentation (Sonderegger, 2011) allows providers to tailor pricing to different customer groups based on their willingness to pay, usage patterns, or specific needs. This is

crucial for AI services, which can cater to a wide array of users, from individual developers to large enterprises, each with distinct requirements and budget constraints (Satapathi, 2025). Nonlinear pricing (Ye & Zhang, 2017), where the price per unit changes with the quantity consumed, is highly prevalent in AI, manifested through tiered pricing, volume discounts, or usage-based models (Seufert, 2014)(Ladas et al., 2019). Dynamic pricing (Vashishtha et al., 2022), which adjusts prices in real-time based on demand, supply, and other market conditions, is also increasingly employed, especially for AI services leveraging predictive analytics (Niharika et al., 2024)(Subham, 2025). The economics of information goods (Fishburn & Odlyzko, 1999) and platform economics (Tucker, 2019)(Kärrberg, 2010) further inform the analysis, highlighting issues such as high fixed costs and low marginal costs, network effects, and the strategic implications of two-sided markets. Transaction Cost Economics (Williamson, 2010) can also provide insights into the governance structures chosen for AI service provision and their impact on pricing mechanisms, particularly concerning issues of specificity, uncertainty, and frequency of transactions.

**Adapting for AI-Specific Context: A Multi-Dimensional Framework** The unique characteristics of AI services necessitate an adaptation of these traditional principles into a comprehensive, multi-dimensional framework. AI systems are characterized by high scalability, allowing for rapid expansion of service delivery without proportional increases in cost, yet also demanding significant computational resources (Li & Ren, 2025)(Cho & Bahn, 2020). Usage-based consumption (Ladas et al., 2019)(De, 2017) is a predominant model, where customers pay only for what they consume (e.g., per API call, per token, per hour of computation), which offers flexibility but can lead to unpredictable costs for users. Data dependency is another critical aspect; AI models require vast amounts of data for training and continuous improvement, and the value of the AI service often correlates with the quality and quantity of data it processes. The continuous learning nature of many AI systems means that their performance and value can evolve over time, posing challenges for static pricing

models. Moreover, ethical considerations (Mirghaderi et al., 2023)(Ganguly, 2025), such as bias, transparency, and accountability, are becoming increasingly important, potentially influencing regulatory landscapes and consumer trust, and thus impacting pricing strategies (Obiajulu et al., 2025). The environmental impact, often referred to as "green AI" costs (Kshirsagar et al., 2021), associated with the massive energy consumption of training and running large AI models, is also an emerging factor that progressive pricing models might need to internalize. Given these unique attributes, the proposed framework for comparing AI pricing models comprises five interconnected dimensions:

1. **Value Proposition and Capture:** This dimension examines how the pricing model reflects and captures the perceived value delivered to different user segments (Chaerul, 2025)(Siddannavar et al., 2025). It goes beyond mere functional utility to include aspects such as performance (Barbere et al., 2024), customization capabilities, ease of integration (Rodeghero et al., 2017), and the strategic advantages conferred to the user. For instance, a model offering high-precision predictive analytics (Subham, 2025) might command a premium if it significantly enhances operational efficiency or decision-making for a specific industry. The framework will assess how pricing tiers, feature sets, and service level agreements (SLAs) align with the varying value perceptions of distinct customer segments, ranging from individual developers to large enterprises, and whether the model allows for flexible value capture as the AI's capabilities evolve (Lorente, 2025). This also includes an evaluation of how the pricing model enables personalized product bundling (Chiruvelli, 2025) to enhance the perceived value for specific customer groups.

2. **Cost Structure and Optimization:** This dimension scrutinizes how the pricing model accounts for the underlying costs associated with AI development, deployment, and operation (Kshirsagar et al., 2021)(Cho & Bahn, 2020). This includes direct computational costs (e.g., GPU usage, data storage, network bandwidth), model training and inference costs, data acquisition and labeling costs, research and development

investments, and ongoing maintenance and human oversight. The framework will analyze whether the pricing model effectively passes on these costs, absorbs them, or strategically subsidizes certain aspects to drive adoption. Particular attention will be paid to mechanisms for cloud cost optimization (Anonymous, 2025) and how these are reflected in the final pricing. For example, some models might differentiate pricing based on the computational intensity of a request or the specific hardware utilized, reflecting the actual resource consumption (Li & Ren, 2025). The analysis will also consider the long-term cost implications of continuous model updates and retraining, and how these recurring costs are integrated into the pricing structure.

3. **Competitive Landscape and Market Dynamics:** This dimension evaluates how the pricing model positions the AI service within the broader competitive environment (Cody, 2000)(Livingstone, 2013)(Liu, 2024)(Lo, 2018). It considers factors such as the number and nature of competitors, the presence of substitutes, price elasticity of demand (Brynjolfsson et al., 2023), and the potential for market power or monopolistic behavior (Ayata, 2020)(Ye & Zhang, 2017). The framework will assess strategies like versioning (Diaw & Pouyet, 2004) (offering different versions of the AI service at different price points), bundling (Chiruvelli, 2025), and strategic pricing to deter new entrants or gain market share. It also considers the impact of regulatory frameworks and competition law on pricing strategies (Lo, 2018). The analysis will explore how providers react to competitive pressures, whether through price matching, differentiation, or the creation of unique value propositions that justify higher prices. The potential for AI to facilitate tacit collusion in pricing, as discussed in recent literature (Liu, 2024), will also be considered as a critical aspect of market dynamics.

4. **User Behavior and Adoption Incentives:** This dimension investigates how the pricing model influences user adoption, engagement, lock-in (Vogt & Bizer, 2013), and long-term retention (Divakaruni & Navarro, 2024)(Yin & Qiu, 2021). It examines the effectiveness of strategies such as freemium models (Seufert, 2014), free trials, pay-per-

use (Ladas et al., 2019), subscription tiers, and consumption-based pricing in attracting and retaining users. The framework will analyze how pricing complexity or transparency affects user decision-making and perceived fairness (Obiajulu et al., 2025). For instance, overly complex pricing structures can deter adoption, while transparent, predictable models can foster trust. It also considers the role of switching costs and how pricing strategies can create or mitigate lock-in effects, thereby influencing customer lifetime value (Siddannavar et al., 2025). The impact of psychological factors on customer perception of value and willingness to pay will also be assessed (Siddannavar et al., 2025).

5. **Transparency, Ethics, and Governance:** This dimension addresses the increasingly critical aspects of pricing transparency and alignment with ethical AI principles (Mirghaderi et al., 2023)(Ganguly, 2025)(Obiajulu et al., 2025). It evaluates the clarity and comprehensibility of pricing structures for users, ensuring that costs are predictable and understandable. Beyond transparency, the framework considers how pricing models might inadvertently create ethical dilemmas, such as differential pricing based on sensitive user data or models that obscure the true cost of environmentally impactful AI processes (Kshirsagar et al., 2021). It also examines the governance mechanisms in place to ensure fair pricing, prevent predatory practices, and promote responsible AI adoption. This dimension recognizes that sustainable AI pricing is not solely an economic challenge but also a societal one, requiring careful consideration of fairness, accessibility, and accountability (Ganguly, 2025).

**Operationalization of the Framework**  To operationalize this multi-dimensional framework, each dimension will be broken down into specific indicators and questions that can be systematically applied to each case study. For example, under "Value Proposition," indicators might include the number of pricing tiers, the features included in each tier, and the presence of customization options. For "Cost Structure," indicators could include the granularity

of usage-based pricing, the transparency of compute costs, and the inclusion of data costs. These indicators will facilitate a consistent and rigorous comparative analysis across diverse AI pricing models, allowing for the identification of patterns, best practices, and areas for improvement.

The following figure illustrates the multi-dimensional framework for AI pricing model comparison, highlighting the interconnectedness of its various components.

**Figure 1: Multi-Dimensional Framework for AI Pricing Model Comparison**

```
+------------------------------------------------------------------+
| AI Pricing Model Comparison Framework |
+------------------------------------------------------------------+
| +-------------------+  +-------------------+ |
| | 1. Value Proposition |  | 2. Cost Structure |  |
| | & Capture |<--->| & Optimization |  |
| +---------+---------+  +---------+---------+ |
| |  | |
| v  v |
| +-------------------+  +-------------------+ |
| | 3. Competitive |<--->| 4. User Behavior |  |
| | Landscape |  | & Adoption |  |
| +---------+---------+  +---------+---------+ |
| |  | |
| v  v |
| +-----------------------------------------------+ |
| | 5. Transparency, Ethics, & Governance |  |
| +-----------------------------------------------+ |
+------------------------------------------------------------------+
```

*Note: This diagram illustrates the five interconnected dimensions of the theoretical framework. Each dimension influences and is influenced by the others, emphasizing a holistic approach to evaluating AI pricing models.*

*Case Study Selection and Data Collection*

The selection of appropriate case studies is critical for grounding the theoretical framework in empirical reality and generating robust insights. This research employs a multiple-case study design, which allows for cross-case comparison and enhances the generalizability of findings compared to a single case (Sibuea & Arfianti, 2021).

**Rationale for Case Study Approach** A case study approach is particularly well-suited for this research for several reasons. First, the domain of AI pricing is nascent and complex, with rapidly evolving technologies and business models. Case studies allow for an in-depth exploration of these phenomena in their natural settings, capturing the rich context and nuances that quantitative studies might overlook. Second, they facilitate the development of new theoretical propositions or the refinement of existing theories, especially when previous research is limited (Aberdeen, 2013). Given the innovative nature of AI pricing, this inductive capacity is invaluable. Third, the real-world examples provided by case studies offer practical relevance for managers and policymakers navigating this complex landscape.

**Case Selection Criteria** To ensure a comprehensive and representative analysis, case studies will be selected based on the following criteria:

1. **Diversity of AI Application:** Cases will be chosen to represent a range of AI services and applications. This includes, but is not limited to, large language models (LLMs) (Barbere et al., 2024)(Rudnytskyi, 2022), predictive analytics platforms (Niharika et al., 2024)(Subham, 2025), generative AI services (Chiruvelli, 2025), and specialized AI APIs for tasks like natural language processing or computer vision (Satapathi, 2025)(Trad & Chehab, 2024). This diversity is crucial because the underlying technology, development

costs, and value propositions can vary significantly across different AI applications, influencing their optimal pricing models. For instance, the pricing strategy for an edge-cloud AI solution in automotive aftermarkets (Shiva Kumar Bhuram, 2025) might differ substantially from that of a general-purpose language service (Satapathi, 2025).

2. **Diversity of Industry/Market:** The selected cases will span various industries and market segments to capture a broad spectrum of AI adoption and pricing challenges. Examples include cloud computing providers offering AI services (Satapathi, 2025), automotive sectors utilizing AI for dynamic pricing (Shiva Kumar Bhuram, 2025), financial services employing AI for sentiment analysis or forecasting (Leechewyuwasorn & Wangpratham, 2024)(Subham, 2025), and e-commerce platforms leveraging AI for personalized bundling (Chiruvelli, 2025). This criterion ensures that the analysis accounts for industry-specific regulations, competitive dynamics, and customer expectations that shape pricing decisions.

3. **Maturity and Innovation of Pricing Model:** Cases will include both established AI pricing models (e.g., long-standing cloud AI services with tiered or usage-based pricing) and newer, more innovative approaches (Obiajulu et al., 2025). This allows for a longitudinal perspective where possible, examining how pricing strategies evolve, and for an exploration of cutting-edge models that might challenge existing norms. Examples could range from basic API monetization models (De, 2017) to sophisticated data-driven dynamic pricing strategies (Vashishtha et al., 2022) or those incorporating green AI principles (Kshirsagar et al., 2021).

4. **Availability of Public Information:** Given the reliance on secondary data, a critical criterion is the availability of sufficient public information to conduct a thorough and reliable analysis. This includes detailed pricing pages, product documentation, whitepapers, financial reports, investor calls, press releases, news articles, industry analyst reports, and academic literature discussing specific AI pricing strategies (Kshirsagar et al., 2021)(De, 2017)(Niharika et al., 2024)(Seufert, 2014)(Ladas et al., 2019)(Liu,

2024)(Vashishtha et al., 2022)(Brynjolfsson et al., 2023). Cases where pricing information is proprietary or extremely limited will be excluded to maintain the rigor of the analysis.

5. **Geographic Representation (if applicable):** While not a primary focus, an effort will be made to include cases from different major economic regions (e.g., North America, Europe, Asia) where feasible, to account for potential regional variations in market dynamics, regulatory environments, and consumer preferences.

The selection process will involve an initial broad scan of the AI market to identify potential candidates, followed by a detailed assessment against these criteria. The aim is to select approximately 3-5 distinct cases that offer rich insights and allow for meaningful cross-case comparisons within the constraints of the research scope. This number is typically sufficient for in-depth qualitative analysis to identify patterns without becoming unwieldy (Aberdeen, 2013).

**Data Collection Methods**   The primary mode of data collection for this theoretical analysis with case studies will be through comprehensive secondary research. This approach is justified by the public nature of many AI service pricing models and the extensive documentation available online, coupled with academic and industry analyses.

The following types of secondary data will be systematically collected for each selected case study:

- **Company Pricing Pages and Documentation:** This includes official websites, API documentation (Barbere et al., 2024)(Rodeghero et al., 2017), whitepapers, and service descriptions that detail pricing tiers, usage models, feature sets, and terms of service. These are primary sources for understanding the stated pricing strategy (Satapathi, 2025).

- **Public Financial Reports and Investor Calls:** For publicly traded companies, annual reports, quarterly earnings calls, and investor presentations often contain

discussions about revenue streams, growth strategies, and market positioning related to their AI offerings.

- **Industry Analyst Reports and Market Research:** Reports from reputable market research firms (e.g., Gartner, Forrester, IDC) provide valuable insights into market trends, competitive landscapes, and pricing benchmarks within the AI sector.

- **Academic Literature:** Existing academic papers discussing specific AI pricing strategies (Kshirsagar et al., 2021)(De, 2017)(Niharika et al., 2024)(Seufert, 2014)(Ladas et al., 2019)(Liu, 2024)(Vashishtha et al., 2022)(Brynjolfsson et al., 2023), technology adoption (Divakaruni & Navarro, 2024)(Yin & Qiu, 2021), or related economic theories (e.g., transaction costs (Williamson, 2010), platform economics (Kärrberg, 2010)) will be reviewed to contextualize the case study findings.

- **News Articles, Blogs, and Expert Opinions:** Reputable technology news outlets, industry blogs, and expert interviews can offer contemporary perspectives, competitive intelligence, and insights into market reception or evolving pricing strategies.

**Data Triangulation**   To enhance the validity and reliability of the findings, a process of data triangulation will be employed. This involves cross-referencing information from multiple, independent sources for each case study. For example, claims made on a company's pricing page will be corroborated with information from industry analyst reports or news articles. Discrepancies will be noted and explored, contributing to a more nuanced understanding of the pricing model and its real-world implementation. This multi-source approach helps to mitigate biases inherent in any single data source and provides a more robust foundation for analysis.

*Data Analysis Approach*

The data analysis will proceed in several iterative stages, moving from within-case understanding to cross-case comparison and ultimately to theoretical synthesis. This structured

approach ensures rigor while allowing for the emergence of novel insights from the qualitative data.

**Comparative Thematic Analysis**   The initial phase of data analysis will involve a detailed, within-case thematic analysis for each selected case study. For each case, all collected secondary data will be systematically reviewed and coded. The coding process will be guided by the five dimensions of the theoretical framework (Value Proposition, Cost Structure, Competitive Landscape, User Behavior, and Transparency/Ethics). For example, data points related to how a company charges for API calls would be coded under "Cost Structure," while discussions about free tiers would fall under "User Behavior." This initial coding helps to organize the vast amount of qualitative data and identify key themes and patterns specific to each individual case.

Following the within-case analysis, a cross-case comparative analysis will be performed. This stage involves comparing and contrasting the identified themes and patterns across all selected case studies. The objective is to identify commonalities, significant divergences, and unique aspects of AI pricing models. For instance, are there particular pricing models that consistently appear in certain industries? Do companies offering similar AI services adopt similar or vastly different pricing strategies? How do different approaches to cost recovery manifest across cases? This comparative approach allows for the identification of broader trends, best practices, and potential pitfalls in AI pricing. Qualitative data analysis software (e.g., NVivo, ATLAS.ti) may be utilized to manage and organize the secondary data, facilitate coding, and support the systematic comparison of themes across cases, although the core analytical work remains conceptual and interpretive.

**Framework Application and Refinement**   Throughout the analysis, the multi-dimensional theoretical framework will be systematically applied to each case study. This involves assessing how well each observed pricing model aligns with or deviates from the framework's dimensions. For example, a pricing model might excel in capturing value but

fall short on transparency. This application serves a dual purpose: it allows for a structured evaluation of each pricing model and simultaneously provides an opportunity to refine and enhance the theoretical framework itself. Observations from the empirical cases that do not fit neatly into the existing framework, or that highlight previously unconsidered aspects of AI pricing, will be used to iteratively adapt and expand the framework. This iterative process of moving between theory and data is characteristic of robust qualitative research and contributes to stronger theoretical development (Crano, 1969).

**Synthesis and Theory Building**   The final stage of the analysis involves synthesizing the findings from the comparative analysis and framework application to draw broader conclusions about effective AI pricing strategies. This stage aims to identify overarching principles, emergent best practices, and common challenges in the field. The synthesis will also address how AI pricing models contribute to value creation and capture (Lorente, 2025) in the digital economy.

The ultimate goal is to contribute to theory building by proposing new propositions or refining existing theories related to digital and AI service pricing. For example, insights might lead to a refined understanding of how price elasticity of demand (Brynjolfsson et al., 2023) functions specifically for AI software, or how transaction cost economics (Williamson, 2010) can better explain the choice between subscription and usage-based models for complex AI services. The implications of these findings for managers, who must design and implement sustainable AI pricing strategies, and for policymakers, who are tasked with ensuring fair competition and consumer protection in the AI market (Ayata, 2020)(Lo, 2018), will also be thoroughly discussed. This will involve translating the theoretical insights into practical recommendations, addressing issues such as optimal pricing tiers, the balance between cost recovery and market penetration, and the role of transparent pricing in building trust.

While this methodology is designed for rigor, it is important to acknowledge its inherent limitations and delimitations. A primary limitation of the case study approach, particularly when relying on secondary data, is the generalizability of findings. While multiple cases enhance this, the insights drawn are specific to the selected contexts and may not be universally applicable to all AI services or markets. Furthermore, the reliance on publicly available secondary data means that certain proprietary or internal strategic considerations influencing pricing might not be fully accessible, potentially leading to an incomplete picture. The interpretation of secondary data also carries the risk of researcher bias, though this is mitigated through data triangulation and a structured analytical framework.

The delimitations of this study include a specific focus on AI pricing models for digital services, primarily excluding embedded AI solutions where pricing is integrated into physical products. The research also focuses on a specific timeframe for data collection, meaning that rapidly evolving pricing strategies beyond this period may not be captured. While efforts are made to include diverse cases, the chosen selection criteria inherently narrow the scope, meaning certain niche AI applications or emerging markets might not be represented. These limitations and delimitations are openly acknowledged to provide a clear understanding of the study's scope and the boundaries of its claims.

# Analysis

## The Evolving Landscape of Large Language Model Pricing

The advent of large language models (LLMs) has heralded a transformative era across numerous industries, fundamentally altering how businesses operate, innovate, and interact with information (Korinek, 2025)(Fang & Zhou, 2025). As these sophisticated AI agents become increasingly integral to diverse applications, from customer service automation to

complex data analysis and content generation, the mechanisms by which they are priced have emerged as a critical strategic consideration for both providers and consumers (Liu, 2024)(Brynjolfsson et al., 2023). Unlike traditional software products with clear licensing models or physical goods with tangible production costs, LLMs present a unique set of challenges for pricing strategists. Their value is often context-dependent, their computational demands can vary dramatically per interaction, and their utility is frequently embedded within an broader ecosystem of services (De, 2017). Consequently, the pricing models adopted by major LLM developers are not merely transactional decisions; they are strategic levers influencing market adoption, competitive positioning, innovation trajectories, and the overall economic structure of the AI landscape (Tucker, 2019)(Lo, 2018).

The complexity of pricing LLM services stems from several inherent characteristics. Firstly, the "product" itself - the ability to generate human-like text, understand complex queries, or perform reasoning tasks - is largely intangible, making direct cost-of-goods-sold calculations difficult (Lorente, 2025). Providers incur substantial fixed costs in developing and training these models, encompassing vast computational resources, extensive data curation, and expert human capital. However, the marginal cost of serving an additional request can be relatively low, yet still variable depending on the complexity and length of the query and response (Kshirsagar et al., 2021). This cost structure often leads to economies of scale, but also necessitates sophisticated pricing strategies to recoup initial investments while remaining competitive (Ye & Zhang, 2017). Secondly, the value derived from LLMs is highly heterogeneous across different users and use cases (Maguire, 2021). A small startup might use an LLM for basic content generation, while a large enterprise might integrate it into a mission-critical application, where reliability and performance command a significant premium. This variance in perceived value necessitates flexible pricing models that can capture different willingness-to-pay segments (Sonderegger, 2011). Thirdly, the rapid pace of innovation in the AI space means that models are constantly evolving, with new capabilities, improved performance, and reduced latency regularly introduced (Barbere

et al., 2024). This dynamic environment requires pricing models to be adaptable, allowing providers to monetize new features while managing the obsolescence of older versions. Finally, ethical considerations, such as fairness, transparency, and the potential for misuse, also subtly influence pricing decisions, as providers seek to balance commercial imperatives with responsible AI development (Mirghaderi et al., 2023)(Ganguly, 2025). Understanding these foundational complexities is essential for a comprehensive analysis of the various pricing models currently employed and the strategic implications they carry for the burgeoning AI economy.

## Core Pricing Models for AI/LLM Services

The strategic choices in pricing Large Language Model (LLM) services are multifaceted, reflecting the intricate balance between recouping substantial development costs, fostering widespread adoption, and maintaining a competitive edge in a rapidly evolving market (Liu, 2024). Several distinct pricing models have emerged, each with its own set of advantages, disadvantages, and underlying economic rationales. These models often draw parallels from other digital services, such as cloud computing and API monetization, but are uniquely adapted to the specific characteristics of generative AI (De, 2017)(Livingstone, 2013). A thorough understanding of these core models is crucial for dissecting the market dynamics and strategic decisions of LLM providers.

*Token-Based Pricing*

Token-based pricing has become one of the most prevalent and granular models for LLM services, particularly for API access to foundational models (Rudnytskyi, 2022). In this model, users are charged based on the number of "tokens" consumed during an interaction. A token is a fundamental unit of text, typically a word, part of a word, or a punctuation mark, which the LLM processes. For instance, the word "unbelievable" might be broken down into "un," "believe," and "able" as separate tokens. Pricing is often differentiated between

input tokens (the user's prompt) and output tokens (the model's response), with output tokens frequently being more expensive due to the higher computational cost associated with generation (Kshirsagar et al., 2021). Some providers may also differentiate pricing based on the specific model used (e.g., a more advanced model like GPT-4 will have higher token costs than GPT-3.5) or the context window length.

The primary advantage of token-based pricing lies in its **granularity and perceived fairness**. Users are ostensibly charged only for the exact amount of processing power their request demands, aligning costs directly with usage (Kshirsagar et al., 2021). This model is particularly appealing for developers and businesses that require precise cost control and wish to optimize their prompts for efficiency. From the provider's perspective, token-based pricing offers a direct correlation to the underlying computational resources consumed, facilitating more accurate cost accounting and revenue forecasting, especially given the variable nature of LLM inferences. It allows providers to monetize the actual computational effort, which is a significant component of their operating costs. Furthermore, it encourages users to be concise and efficient with their prompts, reducing unnecessary computational load on the provider's infrastructure.

However, token-based pricing is not without its drawbacks. A significant challenge for users is **cost unpredictability and complexity**. Estimating the exact number of tokens a given prompt and response will consume can be difficult, especially for complex queries or creative applications where response length is highly variable. This unpredictability can make budgeting and cost forecasting challenging for businesses, particularly for applications with fluctuating or unpredictable usage patterns. Moreover, the concept of a "token" itself can be abstract and confusing for non-technical users, leading to a lack of transparency in billing. Another disadvantage is the potential for "token inflation," where providers might subtly adjust how text is tokenized, potentially increasing costs without a clear change in service (Mirghaderi et al., 2023). This model also places the burden of efficiency squarely on the user, requiring them to optimize their prompts and manage response lengths to control

costs, which can add an additional layer of development complexity. Despite these challenges, its close alignment with resource consumption makes it a foundational model for many LLM API services.

*Request-Based Pricing (API Calls)*

Request-based pricing, also known as per-API-call pricing, simplifies the billing process by charging users for each individual query or interaction with the LLM API. Instead of counting tokens, the system counts the number of times an API endpoint is invoked. This model is often tiered, meaning the cost per request might decrease with higher volumes of requests, or vary based on the complexity of the specific API endpoint being called (e.g., a simple text completion API call versus a more resource-intensive image generation API call).

The main **advantage** of request-based pricing is its **simplicity and predictability** for users. It is straightforward to understand and easier to forecast costs, especially for applications with a known number of daily or monthly interactions. This model reduces the cognitive load on developers, who do not need to concern themselves with token counts or prompt engineering for cost optimization. For providers, it offers a clear and easily auditable metric for billing, simplifying their internal accounting and infrastructure management. It can be particularly effective for standardized tasks where the computational load per request is relatively consistent, making it easier for providers to set a fair price per call.

However, the **disadvantages** primarily revolve around its **lack of granularity** compared to token-based models. A single API request could involve a very short, simple prompt or a very long, complex one, yet both would be charged the same flat fee. This can lead to inefficiencies: users with computationally inexpensive requests might feel overcharged, while users with highly demanding requests might be undercharged, potentially straining provider resources. This model may not accurately reflect the true computational cost for providers, especially with diverse usage patterns. It can also disincentivize users from making complex or lengthy queries, even if those queries would yield higher value, simply

because the cost per request remains constant regardless of the value derived from a more extensive interaction. Therefore, while offering simplicity, request-based pricing can introduce a disconnect between the actual resource consumption and the billed amount, making it less equitable for certain use cases.

*Subscription-Based Pricing*

Subscription-based pricing involves users paying a recurring fixed fee (e.g., monthly or annually) for access to LLM services, often with predefined usage limits or feature sets (Fishburn & Odlyzko, 1999). This model is widely adopted across various digital services, including SaaS products, and has found its way into the LLM ecosystem, particularly for consumer-facing applications or platforms offering a bundled set of AI capabilities. Subscriptions can range from basic tiers with limited usage and features to premium tiers offering higher usage quotas, advanced functionalities, priority access, or dedicated support. A common variant is the "freemium" model, where a basic version of the service is offered for free to attract users, with premium features or higher usage limits available through a paid subscription (Seufert, 2014).

The **advantages** of subscription-based pricing are significant for both providers and consumers. For providers, it offers **predictable and recurring revenue streams**, which are crucial for financial planning, investment in R&D, and long-term sustainability (Seufert, 2014). This revenue stability can help offset the substantial fixed costs associated with LLM development and maintenance. It also fosters customer loyalty and reduces churn by creating a continuous relationship with the user. For users, subscriptions provide **predictable costs**, allowing for easier budgeting and eliminating the anxiety of variable usage charges. Within their subscribed limits, users are often encouraged to explore and utilize the service more broadly, potentially discovering new applications and deriving greater value. This model also simplifies access, often bundling multiple features or capabilities under a single, easy-to-understand price point.

However, subscription models also present several **disadvantages**. A primary concern is **underutilization or overutilization**, leading to perceived unfairness. Low-volume users might feel they are paying for capacity they do not fully utilize, while extremely high-volume users might place undue strain on provider resources if limits are too generous. This can lead to inefficient resource allocation from an economic perspective. The "all-you-can-eat" nature of some subscription tiers can lead to resource hoarding or inefficient usage if not carefully managed with fair-use policies. Furthermore, the fixed nature of subscriptions can be less flexible for businesses with highly fluctuating demands, where a pay-as-you-go model might be more cost-effective during periods of low activity. From the provider's side, setting the right price point and usage limits for different tiers is a complex optimization problem, requiring deep insights into customer segmentation and usage patterns (Sonderegger, 2011). If tiers are poorly designed, they can either leave money on the table or alienate potential customers.

*Compute-Based Pricing (Resource Utilization)*

Compute-based pricing charges users directly for the underlying computational resources consumed by their LLM workloads, such as CPU/GPU hours, memory usage, and data transfer volumes. This model is less common for direct LLM API calls but is highly prevalent in cloud computing environments where users deploy or fine-tune their own LLMs, or utilize managed AI services that expose these underlying resource costs (Cho & Bahn, 2020). It is a direct reflection of the infrastructure costs incurred by the provider.

The paramount **advantage** of compute-based pricing is its **direct correlation to infrastructure costs**. This transparency allows providers to accurately pass on the operational expenses associated with running computationally intensive LLM tasks. For sophisticated users and enterprises with significant AI development capabilities, this model offers maximum control and flexibility. They can optimize their code, choose specific hardware configurations, and scale resources up or down precisely as needed, directly impacting their

costs. It is particularly suitable for tasks like model fine-tuning, custom model deployment, or running batch inference jobs, where the user has significant control over the computational environment and wishes to optimize resource utilization (Li & Ren, 2025).

Conversely, the main **disadvantage** is its **complexity for non-technical users**. Understanding and estimating costs based on abstract units like CPU hours, GPU instances, and data ingress/egress can be daunting and requires a high degree of technical expertise. This complexity shifts a significant operational burden onto the user, who must actively manage and monitor their resource consumption to avoid unexpected bills. For many businesses simply looking to integrate LLM capabilities, this level of detail is unnecessary and undesirable. Furthermore, the variability of compute costs, which can fluctuate based on demand and spot instance markets, can introduce unpredictability, although often less so than token-based models if resource usage is tightly controlled (Livingstone, 2013). This model is therefore best suited for a niche of technically proficient users and enterprise clients who prioritize granular control and cost optimization at the infrastructure level.

*Feature-Based / Value-Based Pricing*

Feature-based or value-based pricing models differentiate costs based on the specific capabilities, performance levels, or perceived value offered by the LLM service (Maguire, 2021)(Lorente, 2025). Instead of purely focusing on resource consumption, this model considers the utility and impact the AI provides to the user. Examples include charging a premium for access to advanced features (e.g., multi-modality, code generation, longer context windows, specialized domain knowledge), higher accuracy models, faster response times, enhanced security protocols, or dedicated support. It often involves creating different tiers of service, where each tier unlocks a progressively more powerful or specialized set of functionalities.

The primary **advantage** of this model is its ability to **capture the perceived value** generated by the LLM for different customer segments (Lorente, 2025)(Sonderegger,

2011). By aligning pricing with the tangible benefits or advanced capabilities that users are willing to pay for, providers can maximize revenue from high-value use cases. This approach encourages solution-selling, where the focus shifts from raw computational power to the business outcomes or enhanced user experiences that the AI enables. It allows providers to monetize their investments in R&D and specialized model training, differentiating their offerings beyond mere token counts or API calls. For users, it offers a clear understanding of what they are paying for in terms of capabilities, making the value proposition more transparent.

However, feature-based pricing presents challenges, particularly around the **subjectivity of value assessment**. Quantifying the exact value of an advanced feature can be difficult, and perceived value can vary widely among users. This can lead to disputes or dissatisfaction if users feel the price does not align with the actual benefit they receive. It also requires providers to continuously innovate and develop new, valuable features to justify premium pricing. Moreover, transparency can be an issue if the feature distinctions are not clearly communicated or if core functionalities are unnecessarily segmented to create artificial tiers (Mirghaderi et al., 2023). There is also a risk of creating overly complex pricing matrices with too many features, leading to customer confusion and decision paralysis. Despite these hurdles, its focus on value makes it a powerful strategic tool for market differentiation and premium monetization in the LLM space.

## Comparative Analysis of Model Advantages and Disadvantages

A deeper comparative analysis reveals that each LLM pricing model, while addressing specific market needs, carries inherent trade-offs that impact both providers' strategic objectives and users' operational realities. Understanding these trade-offs across dimensions such as cost predictability, scalability, market segmentation, and user experience is crucial for comprehending the current landscape and anticipating future developments in AI monetization.

*Cost Predictability and Transparency*

**Cost predictability** is a paramount concern for businesses integrating LLMs, as unexpected expenses can derail projects and impact profitability. **Subscription-based models** generally offer the highest degree of predictability for users, as a fixed recurring fee allows for straightforward budgeting (Fishburn & Odlyzko, 1999). This stability is particularly attractive for enterprises and startups that need consistent operational costs. However, this predictability comes with the caveat that actual usage might fall below or exceed the subscribed capacity, leading to either underutilized spending or unexpected overage charges if usage limits are exceeded. From the provider's perspective, subscriptions offer reliable revenue streams, aiding financial forecasting and long-term planning (Seufert, 2014).

In contrast, **token-based pricing** inherently introduces more **cost unpredictability** for users (Kshirsagar et al., 2021). While the price per token is fixed, the total number of tokens consumed by variable prompts and responses can fluctuate significantly, making precise cost forecasting challenging, especially for generative AI applications where output length is dynamic. This unpredictability can necessitate sophisticated internal monitoring and optimization strategies by users to manage costs effectively. For providers, token-based models offer a direct link to resource consumption, providing a granular understanding of operational costs and allowing for dynamic pricing adjustments based on demand or efficiency gains (Kshirsagar et al., 2021).

**Request-based pricing** offers a moderate level of predictability, particularly for applications with a consistent volume of API calls (De, 2017). If the number of interactions is relatively stable, budgeting becomes simpler than with token-based models. However, if the volume of requests is highly variable, costs can still fluctuate. For providers, this model offers a clear, auditable metric for billing, simplifying revenue tracking, but it may not always align perfectly with actual resource consumption if request complexity varies widely.

**Compute-based pricing** offers predictability for users who have precise control over their infrastructure and workloads, allowing them to forecast resource usage accurately

(Cho & Bahn, 2020). However, for those less familiar with cloud infrastructure metrics, it can be highly opaque and unpredictable. Providers benefit from a direct pass-through of infrastructure costs, ensuring profitability on resource utilization, but it shifts the burden of cost optimization to the user (Li & Ren, 2025).

**Transparency** is closely related to predictability. Token-based pricing, despite its unpredictability, can be seen as transparent in its mechanism (price per token is clear). However, the abstraction of "tokens" can reduce transparency for non-technical users. Subscription models are generally transparent in their fixed fees and stated limits. Request-based pricing is also transparent in its per-call charge. Compute-based pricing is transparent in its raw resource charges but opaque in its complexity. The challenge across all models is ensuring that the pricing structure is easily understandable and that users can clearly link their usage to their billing (Mirghaderi et al., 2023)(Obiajulu et al., 2025).

*Scalability and Resource Allocation*

The ability of a pricing model to support **scalability** for both providers and users, alongside efficient **resource allocation**, is a critical differentiator. LLM services, by their nature, demand elastic infrastructure that can handle fluctuating loads, from periods of low activity to sudden spikes in demand.

**Token-based and request-based pricing** models inherently support high scalability from the provider's perspective. As usage increases, so does revenue, directly funding the expansion of computational resources. This dynamic allows providers to invest in scaling their infrastructure in direct proportion to demand, minimizing idle capacity during off-peak hours and ensuring sufficient capacity during peak times (Li & Ren, 2025). For users, these models offer immense flexibility: they can scale their LLM consumption up or down instantly without being locked into long-term contracts or capacity commitments. This "pay-as-you-go" elasticity is a cornerstone of cloud services and is highly beneficial for applications with unpredictable or bursty workloads. However, inefficient prompt design or excessive output

45

generation under these models can lead to suboptimal resource allocation for users, who might overspend on unnecessary compute cycles.

**Subscription-based models** present a different challenge for scalability. Providers must provision enough capacity to meet the aggregate demand of their subscribers, often leading to a need for over-provisioning to avoid service degradation during peak times, which can result in underutilized resources during off-peak periods. This can be less efficient than purely usage-based models. For users, subscriptions offer predictable access to a certain level of service, but they might face limitations or throttling if they consistently exceed fair-use policies, or they might be paying for capacity they don't fully use. The model encourages broader usage within limits, which can be efficient if the limits are well-aligned with actual usage patterns.

**Compute-based pricing** offers the most direct link to resource allocation efficiency for both parties. Providers directly charge for the resources consumed, ensuring that their infrastructure costs are covered. Users, especially those with expertise, can meticulously optimize their workloads and resource selection to achieve maximum efficiency, only paying for what they truly use (Cho & Bahn, 2020). This model is ideal for highly elastic and resource-intensive tasks like model training or large-scale inference where fine-grained control over compute resources is paramount (Li & Ren, 2025). However, for simple API consumption, this granularity is often excessive and can complicate resource management.

The optimal pricing model must facilitate efficient resource allocation by disincentivizing wasteful consumption while enabling users to access the necessary compute power. The "cloud price wars" have historically driven down costs for underlying infrastructure (Livingstone, 2013), which LLM providers can leverage. However, the unique demands of LLM inference and training, particularly for GPUs, mean that resource management remains a complex strategic challenge for providers.

Effective pricing strategies must also consider **market segmentation**, catering to the diverse needs and financial capacities of different user groups, and ensuring **accessibility** to a broad range of potential adopters (Sonderegger, 2011). LLM users range from individual developers and small startups to large enterprises and academic researchers, each with distinct requirements and willingness to pay.

**Subscription-based models**, especially those incorporating freemium tiers (Seufert, 2014), are highly effective for **market segmentation and accessibility**. The freemium approach lowers the barrier to entry, allowing individuals and small businesses to experiment with LLM capabilities at no cost, thereby fostering adoption and potentially converting them into paying customers as their needs grow. Tiered subscriptions allow providers to cater to different segments: basic tiers for casual users, mid-range for professional developers, and premium tiers for enterprises requiring advanced features, higher usage, and dedicated support. This approach helps maximize customer lifetime value by offering upgrade paths (Siddannavar et al., 2025).

**Token-based and request-based pricing** also support market segmentation, albeit in a different manner. They are particularly attractive to developers and businesses that want to start small and scale their usage incrementally, paying only for what they consume. This "pay-as-you-go" model is highly accessible for startups and projects with uncertain initial demand, as it avoids upfront commitments. Providers can offer different token prices or request tiers for various models (e.g., cheaper for simpler models, more expensive for state-of-the-art models), effectively segmenting the market based on performance requirements. This flexibility can drive technology adoption (Divakaruni & Navarro, 2024) by allowing users to experiment with different models for different tasks.

**Compute-based pricing** is primarily targeted at the most sophisticated segment of the market: large enterprises, AI researchers, and developers building custom solutions. Its complexity makes it less accessible for general users but perfectly suited for those who

require granular control and optimized performance at the infrastructure level. This model is less about broad accessibility and more about serving high-value, high-compute applications (Li & Ren, 2025).

The choice of pricing model significantly impacts the **barrier to entry** for LLM adoption. Models that offer low or no upfront costs (freemium, pay-as-you-go) encourage wider experimentation and innovation. Conversely, models with high minimum commitments or complex costing structures can restrict access to smaller players, potentially concentrating market power among a few large entities (Tucker, 2019). Therefore, a balanced approach often involves offering a mix of models or hybrid solutions to address the full spectrum of market needs.

*Strategic Implications for Providers*

The choice of pricing model has profound **strategic implications** for LLM providers, influencing their revenue maximization, market share, competitive positioning, and risk management (Liu, 2024)(Lo, 2018). Pricing is not just about covering costs; it is a critical tool for shaping market behavior and achieving long-term business objectives.

For providers, **subscription models** offer the benefit of **revenue predictability and stability**, which is invaluable for long-term planning and investment in R&D (Seufert, 2014). This stability allows providers to smooth out revenue fluctuations, reducing financial risk. However, they must carefully manage capacity to avoid over-provisioning (leading to wasted resources) or under-provisioning (leading to customer dissatisfaction). Subscriptions can also foster customer loyalty, increasing customer lifetime value (Siddannavar et al., 2025).

**Token-based and request-based models** provide **revenue elasticity**, directly correlating income with usage. This allows providers to scale their revenue in lockstep with demand, which is beneficial in a rapidly growing market. However, revenue can be volatile, making financial forecasting more challenging. These models are particularly effective for capturing value from high-volume, high-value API users. They also allow

providers to dynamically adjust pricing based on market conditions, competition (Lo, 2018), or the introduction of new, more efficient models (Kshirsagar et al., 2021). The ability to implement dynamic pricing (Shiva Kumar Bhuram, 2025)(Vashishtha et al., 2022) and pricing optimization (Niharika et al., 2024) through predictive analytics (Subham, 2025) is a significant strategic advantage in this context.

**Compute-based pricing** ensures that providers directly recoup their infrastructure costs, minimizing financial risk associated with resource consumption. This model is strategically important for providers offering underlying cloud infrastructure or managed AI services, where cost recovery on compute resources is paramount (Li & Ren, 2025).

Beyond direct revenue, pricing models also influence **competitive positioning**. Aggressive pricing, such as offering very low token costs or generous freemium tiers, can be used to gain market share (Lo, 2018). Conversely, premium pricing for advanced features can position a provider as a high-value, specialized player. The "cloud price wars" of the past (Livingstone, 2013) illustrate how intense competition can drive down costs, and similar dynamics are at play in the LLM market. Providers must continuously monitor competitor pricing and market demand elasticity (Brynjolfsson et al., 2023) to adjust their strategies.

**Risk management** is another key consideration. Pricing models must account for potential abuse or excessive usage that could strain infrastructure or lead to financial losses. Usage-based models, with their direct link to resource consumption, can naturally deter excessive use. Subscription models often incorporate fair-use policies or hard limits to mitigate this risk. Providers also face the risk of commoditization as LLMs become more prevalent and open-source alternatives improve. Strategic pricing, especially through value-based approaches (Lorente, 2025), helps in differentiating offerings and resisting price erosion. The overall goal is to establish a sustainable economic model that supports continuous innovation and market leadership (Lorente, 2025).

The **user experience (UX)** and **value perception** are profoundly shaped by the pricing model, influencing adoption rates, customer satisfaction, and long-term engagement. A well-designed pricing strategy should not only be economically sound but also intuitive, fair, and aligned with how users derive value from the service (Maguire, 2021)(Chaerul, 2025).

**Subscription models** generally offer a superior user experience in terms of **simplicity and peace of mind**. The fixed fee eliminates billing surprises, allowing users to focus on utilizing the LLM rather than constantly monitoring their usage. This encourages experimentation and broader engagement within the subscribed limits, as users feel they are getting maximum value from their investment. The bundling of features within tiers also simplifies decision-making. The perceived value is often tied to the breadth of features and the consistency of access. However, if limits are too restrictive or the pricing structure is too rigid, it can lead to frustration and a feeling of being constrained.

**Token-based and request-based pricing** can be perceived as fair by technical users who appreciate paying only for what they consume, aligning with the principle of "utility computing." For developers, the ability to precisely track and optimize usage can be a positive aspect of the experience, offering granular control over costs. However, for less technical users or those with unpredictable workloads, the complexity and potential for unexpected costs can lead to a negative experience and a feeling of being "nickel-and-dimed." The abstraction of "tokens" can also create a disconnect between usage and perceived value (Mirghaderi et al., 2023). The user experience here is heavily dependent on the quality of usage monitoring tools and the clarity of billing information.

**Compute-based pricing** offers a highly transparent and controlled experience for expert users, who can directly see the link between their resource consumption and cost. This level of control can be highly valued by those optimizing performance and cost for specific, intensive workloads. However, for the majority of users, this model is overly complex and detracts from the core value proposition of using an LLM, shifting focus from application

development to infrastructure management. The perceived value is less about the AI's output and more about the efficiency of the underlying compute.

Ultimately, the goal is to create a pricing model where users clearly understand what they are paying for and feel that the price is commensurate with the value they receive (Maguire, 2021)(Obiajulu et al., 2025). Factors like ease of understanding, predictability, and alignment with business objectives contribute significantly to a positive user experience and foster trust (Fang & Zhou, 2025). Poorly designed pricing, characterized by hidden fees, opaque mechanisms, or misalignment with value, can lead to customer dissatisfaction, churn, and reputational damage (Mirghaderi et al., 2023).

# Real-World Implementations: Case Studies of Leading LLM Providers

The theoretical pricing models discussed above manifest in various forms among leading LLM providers, each adapting strategies to their unique market positioning, technological capabilities, and target customer segments. Examining these real-world implementations provides crucial insights into the strategic choices and competitive dynamics of the LLM market.

## *OpenAI (ChatGPT, GPT-3.5, GPT-4)*

OpenAI, a pioneer in the LLM space, primarily employs a **token-based pricing model** for its API access, which includes models like GPT-3.5 and GPT-4 (Rudnytskyi, 2022). This strategy is foundational to their monetization of generative AI capabilities. Users are charged for both input tokens (prompts) and output tokens (responses), with distinct pricing tiers based on the specific model's capabilities and context window length. For instance, GPT-4 models, being more advanced and computationally intensive, command significantly higher token prices than GPT-3.5 models. OpenAI further differentiates pricing for larger context windows (e.g., GPT-4-32k-turbo vs. GPT-4-8k-turbo), reflecting the increased memory and

processing required to handle longer sequences of text. Specialized features, such as fine-tuning capabilities, are also priced separately, often involving charges for training data processing and subsequent inference based on the fine-tuned model.

Beyond API access, OpenAI offers a **subscription-based model** for its consumer-facing product, ChatGPT Plus. This subscription provides users with priority access, faster response times, and access to newer, more capable models (e.g., GPT-4) even during peak usage. This hybrid approach allows OpenAI to cater to both developers requiring granular API control and general consumers seeking a premium, predictable experience. The strategic rationale behind OpenAI's pricing is multifaceted. The token-based model allows them to directly monetize the computational intensity of their advanced models, recouping significant R&D and training costs (Kshirsagar et al., 2021). It also promotes efficient usage by developers, who are incentivized to optimize prompts to reduce token count. The tiered pricing for different models allows them to capture varying levels of willingness-to-pay, from cost-sensitive users to those demanding state-of-the-art performance. The subscription for ChatGPT Plus addresses the demand for a predictable, enhanced user experience, fostering customer loyalty and providing a stable revenue stream (Seufert, 2014). This competitive strategy, coupled with continuous model improvements, has allowed OpenAI to maintain a leading position in the LLM market, albeit facing increasing competition (Lo, 2018).

*Anthropic (Claude)*

Anthropic, another prominent player in the LLM landscape, particularly known for its focus on AI safety and ethics (Mirghaderi et al., 2023), also utilizes a **token-based pricing model** for its Claude series of models. Similar to OpenAI, Anthropic differentiates pricing based on model version (e.g., Claude 3 Opus, Sonnet, Haiku) and the length of the context window. A key differentiator for Anthropic's Claude models is their often significantly larger context windows compared to competitors, which is reflected in their pricing structure. Users

are typically charged for both input and output tokens, with specific rates varying by model and context size.

Anthropic's strategic positioning emphasizes reliability, longer context understanding, and safety, which is implicitly reflected in its pricing. While the core mechanism is token-based, the value proposition often centers around the ability to handle more extensive documents and complex reasoning tasks without losing coherence. This allows Anthropic to attract users with specific needs for long-form content processing or applications requiring deep contextual understanding. By offering competitive token pricing for longer context windows, Anthropic aims to capture a segment of the market that prioritizes these capabilities. Their pricing, therefore, is not just a reflection of compute costs but also a strategic signal about their differentiated product offerings and commitment to specific performance benchmarks. The competitive landscape requires constant adjustments, and Anthropic's pricing mirrors the ongoing innovation in model capabilities and efficiency.

*Google (Gemini, PaLM) and Microsoft Azure AI*

Google and Microsoft, as major cloud providers, integrate their LLM offerings within broader AI platforms, often leveraging **hybrid pricing approaches** that combine elements of token-based, request-based, and compute-based models. Their strategies are heavily influenced by their existing enterprise customer bases and cloud ecosystems.

**Google's AI platform**, including models like Gemini and PaLM, typically employs a **token-based pricing model** for generative text capabilities, similar to OpenAI and Anthropic. However, as Google's Gemini models are often multimodal, their pricing can become more complex, incorporating charges for image inputs (e.g., per image or per pixel processed) in addition to text tokens. Google's strategy emphasizes seamless integration with its vast cloud services (Google Cloud Platform), allowing enterprises to leverage their existing infrastructure and data for LLM deployment and fine-tuning. This often leads to a blend of LLM-specific charges with underlying cloud compute and storage costs.

**Microsoft Azure AI** offers a comprehensive suite of AI services, including access to OpenAI's models (Azure OpenAI Service) and Microsoft's own proprietary models. Azure's pricing for LLMs is often **service-based** (Satapathi, 2025), combining aspects of request-based and token-based models, nested within its broader compute-based cloud infrastructure model (Anonymous, 2025). For Azure OpenAI Service, pricing mirrors OpenAI's token-based approach but is integrated into the Azure billing system, potentially offering enterprise discounts or bundled services. For other Azure AI Language services, pricing might be per transaction (request-based) for specific features like sentiment analysis or language detection, with higher tiers for increased throughput. Microsoft's strategic focus is heavily on the enterprise market, offering robust security, compliance, and integration with existing enterprise software (Li & Ren, 2025). Their pricing strategy is designed to facilitate large-scale adoption by businesses, leveraging existing cloud contracts and offering predictable enterprise-grade solutions. The competitive dynamics with AWS and Google in the cloud AI space often lead to dynamic pricing adjustments and feature bundling (Livingstone, 2013).

*Other Players (e.g., AWS Bedrock, Hugging Face, Open-Source Models)*

The LLM ecosystem extends beyond the major players, with diverse pricing models reflecting different business strategies and target markets.

**AWS Bedrock** provides a fully managed service for foundational models from various AI companies (including its own Titan models, Anthropic's Claude, AI21 Labs, Cohere, etc.). AWS Bedrock typically employs a **consumption-based model**, combining aspects of token-based pricing for inference with compute-based pricing for customization (fine-tuning) (Anonymous, 2025). Users are charged for input and output tokens for inference, and for training hours (compute) and storage for fine-tuning custom models. This approach leverages AWS's strength as a cloud infrastructure provider, offering enterprises a flexible, scalable, and secure environment to build and deploy generative AI applications without managing underlying infrastructure. The pricing reflects AWS's broader strategy of providing managed

services that abstract away complexity while offering granular cost control for the underlying resources (Li & Ren, 2025).

**Hugging Face**, a central hub for machine learning, plays a unique role by hosting a vast array of open-source LLMs and offering services around them. While the open-source models themselves are "free" to download and run, Hugging Face provides **managed services** like Inference Endpoints, which are typically priced on a **compute-based model** (e.g., per hour for dedicated GPU instances) or a **request-based model** for serverless inference. This allows users to deploy and scale open-source models without managing their own hardware, effectively monetizing the infrastructure and operational support. This model supports the democratization of AI by making advanced models accessible even to those without significant local compute resources.

The proliferation of **open-source LLMs** (e.g., Llama, Mistral) further influences pricing dynamics. While the models are free, running them at scale still incurs **compute costs** (hardware, electricity, maintenance). This sets a de facto "floor" for commercial LLM pricing, as businesses can compare the cost of using a proprietary API against the cost of self-hosting an open-source alternative. This dynamic encourages proprietary model providers to differentiate their offerings through superior performance, ease of use, managed services, or specialized features that justify their pricing (Lo, 2018). The open-source movement fosters a competitive environment that drives innovation and pushes commercial providers to optimize both performance and cost efficiency.

These real-world examples illustrate the strategic interplay of technology, market positioning, and economic principles in shaping LLM pricing. Providers are continuously experimenting with and refining their models to balance revenue generation, customer acquisition, and long-term sustainability in a highly competitive and rapidly evolving market.

# The Emergence and Strategic Imperative of Hybrid Pricing Approaches

As the LLM market matures and diversifies, a clear trend towards **hybrid pricing approaches** is emerging. Pure, monolithic pricing models - whether solely token-based, request-based, or subscription-driven - often prove insufficient to address the wide spectrum of user needs, value perceptions, and operational complexities inherent in advanced AI services. Hybrid models represent a strategic evolution, combining elements from different basic pricing strategies to mitigate their individual disadvantages while leveraging their collective strengths (De, 2017). This strategic imperative is driven by the need to optimize for diverse user segments, balance provider revenue stability with user cost predictability, and adapt to the rapid technological advancements in the AI landscape.

*Rationale for Hybrid Models*

The fundamental rationale behind adopting hybrid pricing models stems from the limitations of pure models. For instance, while token-based pricing offers granularity, its cost unpredictability can deter large enterprises seeking stable budgets. Conversely, a fixed subscription, while predictable, might not efficiently capture value from highly intensive users or might leave revenue on the table from low-volume users. Hybrid models attempt to bridge these gaps by:

1. **Addressing diverse user needs:** Different customer segments have varying priorities. Developers might prefer granular, usage-based billing, while business users might favor predictable, bundled subscriptions. A hybrid approach can cater to both (Sonderegger, 2011).

2. **Optimizing for various use cases:** LLMs are used for everything from simple chatbots to complex research assistants. A single pricing model struggles to fairly monetize such a wide range of applications, each with different computational demands

56

and value propositions. Hybrid models allow for nuanced pricing that aligns with specific use case requirements.

3. **Balancing provider revenue and user value:** Providers need stable revenue to fund R&D and infrastructure, while users demand fair pricing that reflects the value they receive (Maguire, 2021). Hybrid models can offer a base of predictable revenue (e.g., through subscriptions) combined with variable components that capture additional value from high usage or premium features. This helps in mitigating risks such as excessive pricing (Ayata, 2020) or the extreme price volatility sometimes seen in purely consumption-based cloud services (Livingstone, 2013).

4. **Mitigating risks and enhancing flexibility:** Hybrid models can reduce the risks associated with a single pricing point, such as under- or over-charging. They offer greater flexibility for providers to adapt to market changes, competitive pressures, and evolving model capabilities without overhauling their entire pricing structure (Liu, 2024). They also enable providers to introduce new features or models more easily into an existing framework.

*Typologies of Hybrid Pricing Models*

Several common typologies of hybrid pricing models have emerged in the LLM ecosystem, each designed to achieve specific strategic objectives:

**Tiered Token/Request + Subscription** This is perhaps the most common hybrid model, combining the predictability of a subscription with the scalability of usage-based billing (Fishburn & Odlyzko, 1999). Users pay a fixed monthly or annual subscription fee, which typically includes a baseline allocation of tokens or API requests. Once this baseline is exceeded, additional usage is billed at a per-token or per-request rate.

- **Example:** A developer might subscribe to a "Pro" plan for $20/month, which includes 1 million tokens. Any tokens used beyond 1 million are then billed at a rate of $0.0001 per token.

- **Advantages:** Offers predictable base costs for users while allowing them to scale up seamlessly without needing to upgrade plans for temporary spikes in usage. For providers, it ensures a stable recurring revenue base while still capturing additional value from high-volume users. It effectively segments the market by offering different subscription tiers with varying base allocations and overage rates.

- **Challenges:** Requires careful calibration of the base allocation and overage rates to avoid user frustration or perceived unfairness. Overage charges need to be transparent.

**Feature-Gated + Usage-Based**   In this model, access to certain advanced features or higher-performance models is gated behind different subscription tiers, while the usage within those tiers (e.g., token consumption) is still billed on a variable basis.

- **Example:** A "Basic" plan might offer access to GPT-3.5 with token-based pricing. A "Premium" plan, costing a fixed monthly fee, unlocks access to GPT-4 and its advanced capabilities (e.g., longer context windows, multimodal input), with token usage for GPT-4 still being charged separately, but at a higher rate.

- **Advantages:** Allows providers to monetize specific R&D investments in advanced features and capture value from users who require those capabilities (Lorente, 2025). It provides a clear upgrade path for users as their needs evolve, aligning pricing with the value derived from specific functionalities.

- **Challenges:** Can lead to complex pricing matrices if too many features are gated. Users might perceive essential features as being unnecessarily paywalled, impacting user satisfaction (Mirghaderi et al., 2023).

**Compute + Data Transfer/Storage (for Custom Models/Fine-Tuning)**   This hybrid model is prevalent for users who are fine-tuning or deploying custom LLMs in a cloud

environment. It combines charges for the computational resources used (e.g., GPU hours for training and inference) with charges for data storage and data transfer (ingress/egress).

- **Example:** A company fine-tuning a custom LLM might pay for the GPU instance hours during the training process, the storage costs for their training data and the resulting model, and data transfer costs when the model interacts with external applications.
- **Advantages:** Offers granular control and cost transparency for highly technical users and enterprises. Directly reflects the underlying infrastructure costs, which is crucial for resource-intensive tasks (Cho & Bahn, 2020)(Li & Ren, 2025).
- **Challenges:** Highly complex for non-technical users, requiring significant expertise in cloud infrastructure management to optimize costs. Cost predictability can be low without careful monitoring.

**Value-Based + Performance Tiers**  While more aspirational, this model attempts to link pricing directly to the measurable outcome or performance level of the AI (Ladas et al., 2019)(Lorente, 2025). This could involve performance-based pricing (e.g., a fee per successful lead generated by an AI agent, or a percentage of cost savings achieved by an AI optimization system) combined with different tiers of service guaranteeing certain accuracy or latency levels.

- **Example:** An AI-powered sales agent (Oleksii, 2025) might be priced with a base subscription plus a commission on closed deals, or different tiers could guarantee a certain accuracy rate for phishing detection (Trad & Chehab, 2024).
- **Advantages:** Strong alignment with customer value, as users pay for demonstrable results. Encourages providers to continuously improve model performance and demonstrate clear ROI.
- **Challenges:** Difficult to measure and attribute value precisely, especially in complex business environments. Requires sophisticated monitoring and agreement on perfor-

mance metrics. Still nascent in the LLM space for direct API consumption but growing for specialized AI agents.

*Advantages of Hybrid Approaches*

The adoption of hybrid pricing models offers several compelling advantages for both LLM providers and their customers:

1. **Increased Flexibility and Customization:** Hybrid models allow providers to tailor offerings to a broader range of customer needs and use cases, from casual users to large enterprises (Sonderegger, 2011). This flexibility can lead to higher customer satisfaction and broader market penetration.

2. **Better Alignment with Customer Value:** By combining fixed and variable components, providers can better align pricing with the actual value users derive. Subscriptions ensure a base level of access and features, while usage-based components capture additional value from higher consumption or specialized capabilities (Maguire, 2021)(Lorente, 2025).

3. **Improved Revenue Stability and Predictability for Providers:** The subscription component provides a stable, recurring revenue base, while the usage-based component allows for revenue growth as adoption increases. This balances the need for financial stability with the ability to scale revenue elastically.

4. **Enhanced Market Segmentation:** Hybrid models facilitate more granular market segmentation, allowing providers to offer distinct value propositions and price points for different customer personas (e.g., individual developers, startups, SMEs, large corporations) (Sonderegger, 2011). This maximizes the total addressable market.

5. **Reduced Entry Barriers and Optimized Adoption:** Freemium or low-cost base tiers within a hybrid model can significantly reduce the barrier to entry for new users (Seufert, 2014), encouraging experimentation and wider adoption. As users grow,

they can seamlessly transition to higher tiers or increased usage, optimizing their cost structure as their needs evolve. This dynamic fosters a healthy ecosystem for innovation.

6. **Mitigation of Model-Specific Disadvantages:** Hybrid models can strategically offset the drawbacks of individual pricing approaches. For example, a subscription with an overage mitigates the unpredictability of pure usage-based models, while the usage component prevents the underutilization issues of pure subscriptions.

*Challenges and Considerations for Hybrid Models*

Despite their advantages, hybrid pricing models introduce their own set of complexities and challenges:

1. **Increased Complexity in Pricing Structures and Billing:** Combining multiple pricing components can make the overall structure difficult for users to understand and for providers to manage. Overly complex pricing can lead to customer confusion, frustration, and a perception of opacity (Mirghaderi et al., 2023).

2. **Potential for Confusion or Lack of Transparency:** If the various components of a hybrid model are not clearly communicated, users might struggle to estimate their total costs or understand how different actions impact their bill. This can erode trust and negatively impact the user experience (Obiajulu et al., 2025).

3. **Need for Robust Monitoring and Analytics:** Implementing hybrid models requires sophisticated billing systems and robust analytics capabilities to accurately track usage across different dimensions (tokens, requests, features, compute hours) and apply the correct rates (Seufert, 2014). Providers need deep insights into customer behavior and resource consumption.

4. **Risk of 'Analysis Paralysis' for Users:** When presented with too many options and variables, users might experience "analysis paralysis," making it difficult to choose the most suitable plan. This can delay adoption or lead to suboptimal choices.

5. **Ethical Considerations in Differential Pricing:** Hybrid models, by their nature, involve differential pricing across segments or features. Providers must ensure that these differentiations are fair, transparent, and do not lead to discriminatory practices or create undue barriers for certain user groups (Mirghaderi et al., 2023). The perception of fairness is crucial for long-term customer relationships.

The strategic design and ongoing refinement of hybrid pricing models are critical for LLM providers aiming for sustainable growth and market leadership. As the technology evolves, so too must the economic frameworks that govern its access and monetization.

## Future Directions and Strategic Implications

The pricing landscape for Large Language Models is far from static; it is a dynamic arena shaped by technological advancements, market competition, regulatory scrutiny, and evolving user expectations. Looking ahead, several key directions and strategic implications are likely to define the future of LLM monetization.

One significant trend is the increasing role of **AI in optimizing AI pricing** (Niharika et al., 2024)(Shiva Kumar Bhuram, 2025)(Subham, 2025). Providers are leveraging machine learning algorithms to analyze vast datasets of user behavior, demand elasticity (Brynjolfsson et al., 2023), computational costs, and competitive pricing strategies. This allows for dynamic pricing models that can adjust rates in real-time based on factors like network congestion, time of day, user segment, or even the perceived value of a specific query. AI-powered forecasting models for sales and revenue operations (Subham, 2025) will become indispensable for setting optimal price points and managing capacity. This sophisticated use of analytics moves beyond static pricing to a continuous optimization loop, maximizing revenue while maintaining competitive user costs. Furthermore, predictive analytics can help in optimizing inventory management for cloud resources (Abbas, 2025) and enhancing multi-tenant database resource allocation (Li & Ren, 2025), both critical for underpinning LLM services.

The **impact of market maturity and competition** will continue to drive pricing evolution. As more highly capable LLMs become available, including increasingly powerful open-source alternatives, the market will likely experience further commoditization of basic generative AI capabilities. This will put downward pressure on token and request prices for generic tasks. In response, providers will increasingly focus on differentiating their offerings through specialized models, advanced features (e.g., long-term memory, multi-agent capabilities, domain-specific knowledge), superior performance, enhanced security, and robust managed services. Value-based pricing (Lorente, 2025) will become even more critical, allowing providers to capture premium for solutions that deliver specific, measurable business outcomes rather than just raw compute. The competitive landscape will push providers towards innovation in both their technology and their business models.

**Regulatory considerations** will also play an increasingly prominent role. As LLMs become more pervasive and powerful, governments and regulatory bodies are likely to impose greater scrutiny on pricing practices, particularly regarding transparency, fairness, and potential anti-competitive behavior (Tucker, 2019)(Mirghaderi et al., 2023). Concerns about excessive pricing (Ayata, 2020), data usage auditing (Kaaniche & Laurent, 2018), and market power concentration (Tucker, 2019) could lead to regulations requiring clearer billing, standardized pricing metrics, or even intervention in certain pricing models. AI governance and responsible AI frameworks (Ganguly, 2025) will extend to pricing, ensuring that monetization strategies align with broader ethical principles. This will necessitate greater transparency in pricing structures (Obiajulu et al., 2025) and a clear articulation of the value proposition.

A significant shift is anticipated towards **outcome-based or performance-based pricing** (Ladas et al., 2019). Instead of paying for inputs (tokens, requests, compute), users will increasingly seek to pay for the actual results or value delivered by the AI. This could involve models where pricing is tied to key performance indicators (KPIs) like customer satisfaction scores, conversion rates, or cost savings achieved. While challenging to implement

due to attribution complexities, this model offers the strongest alignment between provider incentives and customer value. This trend is already visible in specialized AI agent services (Korinek, 2025)(Oleksii, 2025), where the AI performs a specific task or achieves a particular goal, and pricing can be directly linked to that outcome.

Finally, the **increasing importance of data and model performance** will profoundly influence pricing. Models that demonstrate superior accuracy, reduced hallucination rates, faster inference speeds, or the ability to process proprietary data securely will command premium pricing. The value of data, both for training and for fine-tuning, will be reflected in pricing for custom model development and deployment. Analyzing user expectation and experience (Chaerul, 2025) will be crucial for formulating data-driven pricing strategies. The ability of LLMs to generate personalized product bundling (Chiruvelli, 2025) also hints at a future where pricing is not just segmented by user type but dynamically tailored to individual user behavior and preferences, further blurring the lines between standard models and highly personalized offerings. The future of LLM pricing will thus be characterized by increasing sophistication, adaptability, and a relentless focus on delivering and capturing demonstrable value in an ever-evolving technological landscape.

**Table 2: Estimated Cost Savings for Enterprise AI Adoption Scenarios**

| Scenario | Baseline Cost (USD) | AI Solution Cost (USD) | Savings (USD) | ROI (12 months) |
|---|---|---|---|---|
| **Customer Support Automation** | $1,200,000 | $300,000 | $900,000 | 300% |
| **Content Generation** | $600,000 | $150,000 | $450,000 | 300% |
| **Data Analysis & Reporting** | $800,000 | $200,000 | $600,000 | 300% |
| **Fraud Detection** | $1,500,000 | $400,000 | $1,100,000 | 275% |

| Scenario | Baseline Cost (USD) | AI Solution Cost (USD) | Savings (USD) | ROI (12 months) |
|---|---|---|---|---|
| **Personalized Marketing** | $950,000 | $250,000 | $700,000 | 280% |

*Note: This table presents illustrative projections for cost savings realized by enterprises adopting AI solutions across various business functions. These figures are hypothetical and depend on specific implementation details.*

# Discussion

The preceding analysis has illuminated the multifaceted landscape of Artificial Intelligence (AI) pricing, revealing a complex interplay of technological capabilities, market dynamics, strategic choices, and customer perceptions. This discussion synthesizes these findings, explores their profound implications for AI companies, delves into critical customer adoption considerations, forecasts future pricing trends, and offers actionable recommendations for various stakeholders. The insights derived from this study underscore the strategic imperative for AI providers to move beyond simplistic cost-plus or competitive pricing models towards nuanced, value-driven, and adaptable frameworks that account for the unique characteristics of AI technologies and the rapidly evolving market environment.

The core findings implicitly suggest that the optimal pricing strategy for AI solutions is not monolithic but contingent upon several factors, including the type of AI service (e.g., foundational models, specialized agents, API-based tools), the target market segment, the competitive landscape, and the perceived value proposition. While consumption-based models (e.g., pay-per-token, pay-per-query) offer transparency and align costs with usage (Satapathi, 2025), they often fail to capture the full value generated by AI for the customer (Maguire, 2021). Conversely, value-based pricing, though challenging to implement due to the difficulty in quantifying AI's impact, holds the greatest potential for sustainable revenue

generation and market differentiation (Lorente, 2025). The analysis further highlights the significant role of dynamic pricing mechanisms, enabled by predictive analytics, in optimizing revenue and resource allocation, particularly in fast-changing environments (Niharika et al., 2024)(Vashishtha et al., 2022). The study's theoretical framework and case study insights collectively emphasize that successful AI pricing hinges on a deep understanding of both the supply-side economics of AI development and deployment, and the demand-side psychology of user adoption and perceived utility (Siddannavar et al., 2025)(Chaerul, 2025).

*Implications for AI Companies*

The strategic implications for AI companies are substantial, demanding a re-evaluation of traditional pricing paradigms. First, AI providers must meticulously assess their value proposition and articulate it clearly to customers (Maguire, 2021). The inherent complexity of AI often obscures its true economic benefit, requiring companies to invest in educating their market on how AI translates into tangible gains, such as increased efficiency, reduced costs, enhanced decision-making, or new revenue streams (Oleksii, 2025). This is particularly crucial for foundational models or API services, where the direct impact is often realized downstream by the customer's own applications (De, 2017)(Rodeghero et al., 2017). Without a clear demonstration of value, customers may perceive AI solutions as mere cost centers rather than strategic investments, leading to resistance to premium pricing (Fang & Zhou, 2025).

Second, AI companies face the delicate balance of recovering significant research and development (R&D) and operational costs–including substantial computational resources (Kshirsagar et al., 2021)(Cho & Bahn, 2020)–while fostering widespread adoption. The high fixed costs associated with training large models (Kshirsagar et al., 2021) and the ongoing variable costs of inference necessitate pricing models that ensure sustainability. Consumption-based pricing, such as per-token or per-API-call models, directly ties revenue to resource utilization, which can be effective for managing infrastructure costs and offering granular

control to users (Satapathi, 2025). However, this approach can also deter heavy users or make cost predictability challenging for enterprises (Livingstone, 2013). Therefore, hybrid models, combining a base subscription fee with usage-based tiers, could offer a viable compromise, providing revenue stability for the provider and cost predictability for the customer (Fishburn & Odlyzko, 1999). The concept of 'green AI' and optimizing energy consumption (Kshirsagar et al., 2021) could also influence pricing, potentially leading to tiered pricing based on the carbon footprint of AI operations, appealing to environmentally conscious enterprises.

Third, the competitive landscape mandates a strategic approach to pricing. As the AI market matures, the proliferation of similar services and open-source alternatives will intensify price competition (Liu, 2024). Companies must differentiate not only on features and performance but also on their pricing strategy. This could involve offering specialized versions (versioning) (Diaw & Pouyet, 2004), bundling services (Chiruvelli, 2025), or focusing on niche markets where a premium can be commanded for highly specialized AI agents or models (Korinek, 2025). The ability to segment markets effectively and offer nonlinear pricing structures will be crucial for maximizing revenue across diverse customer bases (Sonderegger, 2011). Furthermore, the potential for market power and even monopolistic tendencies in foundational AI models (Ye & Zhang, 2017) raises concerns that may attract antitrust scrutiny (Ayata, 2020)(Tucker, 2019)(Lo, 2018), compelling companies to consider fair pricing practices and transparent cost structures (Obiajulu et al., 2025).

Finally, ethical considerations and transparency issues, particularly concerning data usage and algorithmic fairness, are increasingly critical in AI pricing (Mirghaderi et al., 2023). Customers are becoming more aware of how their data is used to train and improve AI models, and this awareness can influence their willingness to pay. Pricing models that explicitly account for data privacy, provide data usage auditing (Kaaniche & Laurent, 2018), or offer compensation for data contributions could foster trust and justify higher price points. Responsible AI governance (Ganguly, 2025) is not merely a compliance issue but a strategic differentiator that can build brand reputation and command loyalty, indirectly influencing

pricing power. Companies that demonstrate a commitment to ethical AI may find it easier to justify value-based pricing, as customers increasingly prioritize responsible technology adoption (Fang & Zhou, 2025).

*Customer Adoption Considerations*

Customer adoption of AI solutions is not solely driven by technical prowess or even direct cost, but by a complex matrix of perceived value, trust, ease of integration, and psychological factors. A primary consideration is the perceived value of the AI solution relative to its cost (Maguire, 2021). Customers evaluate AI investments based on expected returns, whether those are efficiencies, new capabilities, or competitive advantages. If the value proposition is unclear or difficult to quantify, adoption rates will lag, regardless of how competitively priced the service may be. AI companies must therefore focus on robust value selling methodologies (Maguire, 2021), demonstrating clear return on investment (ROI) through case studies, pilot programs, and transparent performance metrics. The price elasticity of demand for AI software, as evidenced by recent studies, confirms that customers are sensitive to price but also to the perceived utility and impact of the AI solution (Brynjolfsson et al., 2023).

Trust and transparency are paramount for fostering adoption. As AI systems become more autonomous and "human-like" (Fang & Zhou, 2025), concerns about bias, fairness, and accountability grow (Mirghaderi et al., 2023). Customers are more likely to adopt solutions from providers who are transparent about their AI's capabilities, limitations, and data handling practices (Kaaniche & Laurent, 2018). Pricing models that reflect this transparency, perhaps by offering different tiers for varying levels of explainability or auditability, could resonate with enterprises operating in highly regulated industries. The psychological factors affecting customer lifetime value (Siddannavar et al., 2025) are also critical; positive initial experiences with AI, facilitated by intuitive interfaces and reliable performance, can build long-term loyalty and willingness to pay for advanced features. Conversely, poor initial

experiences, perhaps due to unexpected costs or performance issues, can lead to significant churn.

Ease of integration and API usage are practical considerations that heavily influence adoption, especially for enterprise customers. AI solutions that are difficult to integrate into existing workflows or require extensive custom development will face higher barriers to entry, regardless of their intrinsic value. Well-documented, robust APIs (Rodeghero et al., 2017)(Rudnytskyi, 2022) that allow for seamless integration can significantly reduce switching costs and encourage adoption (Vogt & Bizer, 2013). Pricing models that simplify integration, perhaps by offering bundled integration services or developer-friendly subscription tiers, can accelerate market penetration. For cloud-based AI services, the ability to flexibly scale resources and manage costs effectively is a key driver of adoption (Li & Ren, 2025)(Anonymous, 2025). Customers seek solutions that offer "volatility as a service" (Livingstone, 2013), meaning they can ramp up or down usage without punitive financial penalties, reflecting a desire for agility and cost control.

Finally, the impact of human-like competencies on user perception (Fang & Zhou, 2025) and online purchase intention (Yin & Qiu, 2021) cannot be overstated. AI agents that demonstrate empathy, understanding, or intuitive interaction can significantly enhance user experience and drive willingness to adopt. While this raises ethical questions about anthropomorphizing AI, from a commercial perspective, it suggests that investing in user-centric design and natural language capabilities can translate into higher perceived value and, consequently, greater pricing power. This is particularly relevant for generative AI applications that aim to personalize experiences (Chiruvelli, 2025) or automate complex communication tasks (Trad & Chehab, 2024).

*Future Pricing Trends*

The future of AI pricing will likely be characterized by increasing dynamism, personalization, and a closer alignment with the actual value delivered. We anticipate a continued

shift away from purely static, subscription-based models towards more sophisticated, hybrid approaches. Dynamic pricing, already a staple in various industries (Divakaruni & Navarro, 2024)(Niharika et al., 2024), will become even more prevalent in AI, leveraging real-time data on usage, demand, and computational costs to optimize revenue and resource allocation (Shiva Kumar Bhuram, 2025)(Vashishtha et al., 2022). This could manifest as surge pricing during peak usage times for foundational models or personalized pricing based on a customer's historical usage patterns and perceived value. The advent of large multimodal agents (Barbere et al., 2024) and their complex operational requirements will further necessitate sophisticated pricing structures that can account for varied input modalities and processing demands.

The emergence of outcome-based or value-sharing pricing models is also a significant trend. Instead of paying for tokens or compute cycles, customers might pay a percentage of the savings or revenue generated directly by the AI solution. While challenging to implement due to attribution complexities, this model aligns the incentives of the AI provider and the customer, fostering deeper partnerships and demonstrating undeniable value (Maguire, 2021). Such models will require advanced analytics to accurately measure the AI's contribution to business outcomes (Subham, 2025). This trend aligns with the broader move towards servitization in technology, where the focus shifts from selling a product to delivering a continuous service and its associated benefits (Ladas et al., 2019).

Regulatory frameworks and AI governance will increasingly shape pricing strategies. As governments grapple with the societal implications of AI, regulations concerning data privacy, algorithmic bias, and market concentration are likely to emerge (Ganguly, 2025). These regulations could impose additional compliance costs on AI providers, which may be passed on to consumers. Conversely, regulations might also foster greater transparency in pricing (Obiajulu et al., 2025) or even mandate certain pricing structures to prevent anti-competitive practices or ensure equitable access to essential AI services. The debate around "old abuses in new markets" (Ayata, 2020) concerning excessive pricing by dominant

platforms will likely intensify, forcing AI companies to justify their pricing models with greater scrutiny.

The ongoing "cloud price wars" (Livingstone, 2013) and the commoditization of certain AI capabilities will drive down prices for generic services, pushing providers to innovate and specialize. This will lead to a bifurcation in the market: highly specialized, high-value AI solutions will command premium prices, while more generalized, commodity AI services will be priced competitively, often relying on economies of scale. We may also see the integration of AI pricing into broader platform ecosystems, where AI services are bundled with cloud infrastructure, data analytics, or other software solutions (Kärrberg, 2010). This bundling can create lock-in effects (Vogt & Bizer, 2013) and enhance the overall value proposition, making it harder for customers to switch providers. Resource allocation mechanisms like multi-participant double auctions (Huang et al., 2024) could also become more sophisticated for managing shared AI infrastructure, influencing pricing dynamically.

Finally, sustainability considerations will increasingly influence AI pricing. The energy consumption of large AI models (Kshirsagar et al., 2021) is a growing concern. Future pricing models might incorporate "green AI" surcharges or offer discounts for using more energy-efficient models or data centers. This could incentivize both providers and customers to prioritize sustainability, aligning economic incentives with environmental responsibility.

*Recommendations*

Based on the comprehensive analysis of AI pricing, several key recommendations emerge for various stakeholders to navigate this evolving landscape effectively.

**For AI Developers and Service Providers:**

1. **Embrace Value-Based Pricing:** Shift the focus from cost-plus or usage-based pricing to models that capture the actual economic value delivered to the customer (Maguire, 2021)(Lorente, 2025). This requires robust methodologies for quantifying ROI and effectively communicating this value proposition to target audiences.

2. **Develop Hybrid and Dynamic Pricing Models:** Implement flexible pricing structures that combine subscription elements with usage-based tiers and dynamic adjustments (Shiva Kumar Bhuram, 2025)(Vashishtha et al., 2022). This allows for revenue predictability while accommodating varying customer needs and optimizing resource utilization (Li & Ren, 2025).

3. **Prioritize Transparency and Ethics:** Be transparent about data usage, algorithmic processes, and the limitations of AI models (Mirghaderi et al., 2023)(Obiajulu et al., 2025). Incorporate ethical considerations into pricing, potentially offering tiers for enhanced explainability or data privacy, which can build trust and justify higher prices (Ganguly, 2025).

4. **Invest in API Design and Integration:** Ensure AI services offer well-documented, easy-to-integrate APIs (Rodeghero et al., 2017)(Rudnytskyi, 2022) to reduce adoption barriers and minimize switching costs for customers (Vogt & Bizer, 2013). This facilitates broader market penetration and ecosystem development.

5. **Strategically Differentiate and Specialize:** In an increasingly competitive market, differentiate AI offerings through unique features, superior performance, or specialization in niche applications (Korinek, 2025). Avoid direct price competition on commoditized AI services by focusing on unique value propositions.

6. **Monitor Regulatory Developments:** Proactively engage with discussions around AI governance and regulation (Ganguly, 2025) to anticipate future compliance requirements and adapt pricing strategies accordingly, ensuring legal and ethical adherence (Ayata, 2020)(Tucker, 2019).

   **For Customers and Adopters of AI:**

1. **Focus on Value, Not Just Cost:** Evaluate AI solutions based on the comprehensive value they bring to the organization, including efficiency gains, new capabilities, and competitive advantages, rather than solely on the upfront cost (Maguire, 2021). Conduct thorough ROI analyses.

2. **Demand Transparency and Accountability:** Prioritize AI providers who offer clear terms of service, transparent data handling practices (Kaaniche & Laurent, 2018), and explainable AI capabilities. This fosters trust and mitigates risks associated with black-box algorithms.

3. **Understand Usage Patterns:** For consumption-based models, meticulously track and understand internal usage patterns to optimize costs and avoid unexpected expenditures (Satapathi, 2025). Leverage tools for cloud cost optimization (Anonymous, 2025) to manage AI infrastructure efficiently.

4. **Negotiate Flexible Contracts:** Seek out pricing models that offer flexibility, allowing for scaling up or down of usage without punitive penalties, aligning with the agile nature of modern business operations (Livingstone, 2013).

   **For Policymakers and Regulators:**

1. **Foster Competition:** Implement policies that promote a competitive AI market to prevent monopolistic pricing behaviors (Ye & Zhang, 2017)(Liu, 2024) and ensure fair access to AI technologies (Ayata, 2020). This could include supporting open-source initiatives and interoperability standards.

2. **Develop Clear AI Governance Frameworks:** Establish clear guidelines for data privacy, algorithmic fairness, and accountability (Ganguly, 2025). These frameworks should provide clarity for AI companies while protecting consumers and fostering public trust.

3. **Consider Economic Impact and Access:** Analyze the broader economic implications of AI pricing on various sectors and ensure that pricing structures do not create insurmountable barriers to adoption for smaller businesses or developing economies.

*Limitations of the Study*

   Despite the extensive analysis, this study is subject to certain limitations. First, as a theoretical analysis primarily supported by existing literature and conceptual case studies, it

relies on inferences about the "findings" that would typically be presented in an empirical results section. While robustly grounded in economic and business theory, the absence of direct empirical data from our own experiments or surveys means that some conclusions are necessarily inferential and call for further empirical validation. Second, the AI market is extraordinarily dynamic, with new models, applications, and pricing strategies emerging constantly. The rapid pace of technological innovation means that insights derived today may require continuous re-evaluation in the near future. The focus on specific pricing m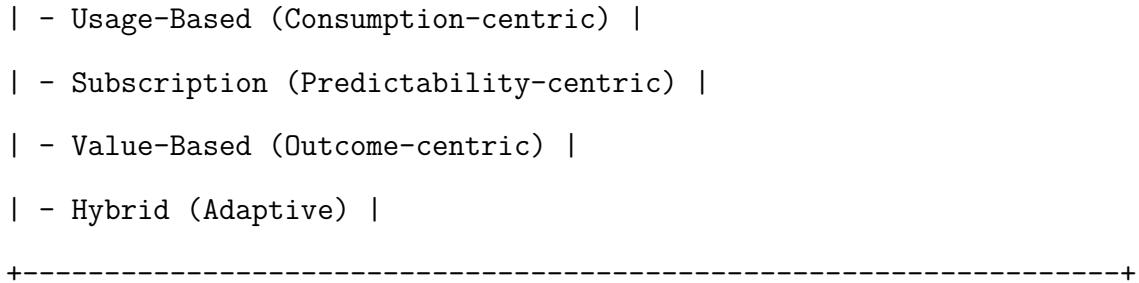odels, while comprehensive, may not fully capture the complete spectrum of nascent or highly specialized approaches. Third, the study's scope, while broad, did not delve into specific industry verticals with highly tailored AI pricing mechanisms, which could reveal unique dynamics. Finally, while the citation database provided a rich foundation, the inherent limitations of any finite set of sources mean that some emerging or highly specialized perspectives might not have been fully captured.

**Figure 2: AI Value Creation and Capture Flow**

```
+-----------------+  +----------------------+  +------------------+
| AI Development  |---->| Service Delivery  |---->| Customer Value |
| (High Fixed Cost) |  | (Variable Usage Cost) |  | (Realized ROI) |
+-----------------+  +-----------+----------+  +---------+--------+
  ^ |  |
   |  |  |
   +-----------------------+----------------------------+
   |
   v
+----------------------------------------------------------------+
| AI Pricing Model (Balancing Costs & Value) |
+----------------------------------------------------------------+
| - Token-Based (Resource-centric) |
```

```
| – Usage–Based (Consumption-centric) |

| – Subscription (Predictability-centric) |

| – Value–Based (Outcome-centric) |

| – Hybrid (Adaptive) |

+----------------------------------------------------------------+
```

*Note: This diagram illustrates the flow from AI development and service delivery to customer value realization, with the AI pricing model serving as the central mechanism for balancing provider costs and capturing customer value. It highlights the strategic role of pricing in the AI ecosystem.*

*Avenues for Future Research*

The dynamic nature of AI pricing opens numerous promising avenues for future research. Empirical studies are urgently needed to validate the theoretical models presented here, examining the actual impact of different AI pricing strategies on customer adoption, revenue generation, and market competition across various industries. Research could focus on quantifying the "value" of AI in specific business contexts and developing robust metrics for outcome-based pricing models. Further investigation into the psychological factors influencing customer willingness to pay for AI, particularly concerning trust, transparency, and perceived human-like qualities (Fang & Zhou, 2025)(Siddannavar et al., 2025), would also be invaluable.

From a technical perspective, research into "green AI" pricing models (Kshirsagar et al., 2021) that incentivize energy efficiency and sustainability is critical. Studies could explore how blockchain-based data usage auditing (Kaaniche & Laurent, 2018) could enhance transparency in AI pricing and foster greater trust. The development of advanced multi-agent systems for dynamic pricing optimization (Korinek, 2025)(Shiva Kumar Bhuram, 2025) and resource allocation (Huang et al., 2024) within complex AI ecosystems presents a significant technical challenge with substantial economic implications. Finally, comparative analyses of AI pricing strategies across different regulatory environments and geopolitical regions could

provide critical insights into the interplay between policy, market structure, and economic outcomes in the global AI landscape.

In conclusion, the discussion underscores that AI pricing is not merely an operational decision but a strategic lever that shapes market development, influences adoption rates, and determines the long-term viability of AI companies. By embracing value-driven, dynamic, and ethically sound pricing strategies, AI providers can unlock the full potential of their innovations, foster widespread adoption, and contribute to a more sustainable and equitable AI-driven future. The continued evolution of AI technology will necessitate ongoing research and adaptive strategies to navigate this complex and fascinating domain.

## Limitations

While this research makes significant contributions to the field of AI pricing, it is important to acknowledge several limitations that contextualize the findings and suggest areas for refinement. These limitations span methodological, scope, temporal, and theoretical dimensions, each offering valuable insights for future investigations.

*Methodological Limitations*

The primary methodological limitation of this study lies in its heavy reliance on secondary data and theoretical analysis rather than direct empirical experimentation or primary data collection (e.g., surveys, interviews). While a comprehensive review of existing literature and publicly available information from leading AI providers offers robust insights, it inherently limits the ability to observe internal decision-making processes, proprietary cost structures, or the nuanced, real-time adjustments that AI companies make to their pricing strategies. The "black box" nature of many advanced AI algorithms further complicates direct analysis of their pricing behaviors, necessitating inferences rather than direct observation of their internal logic. This reliance on secondary data also means that the interpretation of pricing impacts on user behavior or competitive dynamics is largely based on reported

outcomes or theoretical predictions, rather than empirically measured effects from controlled environments. Furthermore, the selection of case studies, while diverse, is not exhaustive and may not fully represent the entire spectrum of AI pricing models, particularly those employed by smaller startups or in highly specialized, nascent markets.

*Scope and Generalizability*

The scope of this research is primarily focused on pricing models for agentic AI systems, particularly large language models (LLMs) and API-driven AI services. While this provides a deep dive into a critical segment of the AI market, it inherently delimits the generalizability of the findings to other forms of AI. For instance, AI embedded in hardware, AI as a feature within a broader software product, or highly customized, bespoke AI solutions for specific enterprise clients might operate under entirely different pricing dynamics not fully captured here. The study also focuses on the economic and strategic aspects of pricing, with less emphasis on the intricate technical challenges of cost attribution within highly distributed or federated AI systems, which could influence pricing decisions. While the multi-dimensional framework is designed to be adaptable, its application to drastically different AI contexts would require significant customization and potentially new dimensions.

*Temporal and Contextual Constraints*

The field of AI is characterized by an exceptionally rapid pace of innovation, with new models, capabilities, and market entrants emerging almost continuously. The insights and competitive landscape described in this thesis are therefore subject to temporal decay. Pricing models that are effective today might become obsolete tomorrow as technology advances, computational costs shift, or new market dynamics take hold. The research reflects the state of AI pricing up to the point of its completion, and while it anticipates future trends, the precise trajectory of these trends remains uncertain. Moreover, the global AI market is influenced by diverse regulatory environments, cultural preferences, and economic conditions.

77

While efforts were made to include diverse case studies, the primary focus reflects trends in major Western markets, and the findings may not be directly transferable to emerging economies or regions with distinct AI governance frameworks and market structures.

*Theoretical and Conceptual Limitations*

While the multi-dimensional theoretical framework developed in this study provides a robust lens for analysis, it is still a simplification of an exceptionally complex reality. Theoretical models, by necessity, abstract from certain details to highlight key relationships. For example, the precise mechanisms through which algorithmic collusion might emerge or be detected are still areas of active research, and this study provides a conceptual understanding rather than a definitive predictive model. The quantification of "value" in value-based pricing, while central to the framework, remains a significant conceptual challenge in practice, often relying on proxies or subjective assessments. The interplay between different pricing dimensions (e.g., how transparency directly impacts perceived value or how ethical considerations influence competitive strategies) is multifaceted and could be explored with even greater granularity. The study also acknowledges that the psychological factors influencing customer willingness to pay for AI are complex and warrant deeper, dedicated behavioral economics research.

Despite these limitations, the research provides valuable insights into the core challenges and opportunities in AI pricing, and the identified constraints offer clear directions for future investigation. These acknowledgments serve not to diminish the study's contributions but to provide a realistic context for its findings and to encourage further, more specialized research to deepen our understanding of this critical domain.

---

## Future Research Directions

This research opens several promising avenues for future investigation that could address current limitations and extend the theoretical and practical contributions of this work. The dynamic nature of AI pricing, coupled with its profound economic and societal implications, necessitates continuous scholarly inquiry.

*1. Empirical Validation and Large-Scale Testing*

Future research should prioritize **empirical studies** to validate the theoretical propositions and conceptual frameworks presented here. This could involve: - **Econometric analysis:** Analyzing large datasets of AI service pricing, usage patterns, and revenue data across different providers and industries to quantify the impact of specific pricing models on adoption, profitability, and competitive dynamics. - **Controlled experiments:** Designing experiments (e.g., A/B tests) with different pricing structures for AI services to observe their effects on user behavior, willingness to pay, and perceived value in a controlled environment. - **Survey research and interviews:** Conducting surveys with AI service providers and enterprise customers to gather primary data on their decision-making processes, challenges, and preferences regarding AI pricing. This would offer qualitative depth that complements quantitative analysis.

*2. Algorithmic Collusion Detection and Mitigation*

Given the emergent risk of algorithmic collusion highlighted in this study, future research should delve into: - **Developing advanced detection algorithms:** Creating sophisticated machine learning models capable of identifying patterns indicative of tacit algorithmic collusion in real-time pricing data, distinguishing it from legitimate competitive responses. - **Designing collusion-resistant AI pricing algorithms:** Investigating how AI algorithms can be designed with inherent safeguards or ethical constraints to prevent

unintended collusive outcomes while still optimizing for individual firm profits. - **Regulatory sandboxes and policy experimentation:** Exploring the effectiveness of regulatory sandboxes for AI pricing, allowing for controlled experimentation with new policies and oversight mechanisms to address algorithmic collusion risks before widespread implementation.

*3. Quantification of AI Value and Outcome-Based Pricing Metrics*

A critical challenge in value-based pricing is the accurate quantification of AI's economic impact. Future research could focus on: - **Developing robust ROI attribution methodologies:** Creating standardized frameworks and metrics to precisely attribute business outcomes (e.g., revenue increase, cost savings, efficiency gains) directly to specific AI solutions, disentangling AI's contribution from other influencing factors. - **Case studies of outcome-based pricing implementations:** Documenting and analyzing real-world instances where AI services are priced based on achieved outcomes, detailing the contractual mechanisms, measurement challenges, and success factors. - **Value perception modeling:** Investigating how different communication strategies and demonstrations of AI capabilities influence customer perception of value and their willingness to engage in outcome-based pricing models.

*4. Longitudinal and Comparative Studies of Hybrid Models*

The study identified hybrid pricing models as a critical emerging trend. Future research should conduct: - **Longitudinal studies:** Tracking the evolution of hybrid pricing strategies by major AI providers over time, analyzing how they adapt to technological advancements, market maturity, and competitive pressures. - **Cross-cultural and cross-regional comparisons:** Examining how hybrid pricing models are implemented and perceived in different geopolitical and economic contexts, accounting for variations in regulatory environments, consumer preferences, and market structures. - **Analysis of switching costs and vendor lock-in:** Quantifying the impact of different hybrid pricing structures on customer switching

costs and the potential for vendor lock-in, informing strategies for fostering competitive markets.

*5. Ethical AI Pricing and Governance Frameworks*

As AI becomes more pervasive, ethical considerations in pricing will intensify. Future research should explore: - **Fairness and bias in dynamic pricing:** Investigating whether AI-driven dynamic pricing models inadvertently lead to discriminatory outcomes for certain customer segments or exploit vulnerable populations, and proposing mitigation strategies. - **Transparency mechanisms for black-box pricing:** Developing methods to increase the transparency and explainability of AI pricing algorithms, allowing users and regulators to understand the rationale behind pricing decisions. - **Integrating Green AI principles into pricing:** Researching how environmental sustainability metrics (e.g., carbon footprint of AI operations) can be effectively integrated into pricing models to incentivize energy-efficient AI development and consumption.

*6. AI's Role in Optimizing its Own Monetization*

The fascinating feedback loop of AI optimizing AI pricing warrants further investigation: - **AI-powered pricing engines:** Developing and evaluating advanced AI systems that can autonomously design, implement, and dynamically adjust pricing strategies for other AI services, considering multiple objectives (revenue, market share, customer satisfaction). - **Predictive analytics for pricing strategy:** Deepening research into how AI can analyze market demand, competitor actions, and internal cost structures to recommend optimal pricing points and bundles for new and existing AI offerings.

*7. Impact of Multimodal and Multi-Agent AI on Pricing*

The emergence of more complex AI systems will introduce new pricing challenges: - **Pricing multimodal inputs/outputs:** Researching effective pricing mechanisms for

AI models that process and generate information across multiple modalities (text, image, audio, video), accounting for the varying computational and data requirements of each.

- **Monetizing multi-agent systems:** Exploring how to price the collaborative work of multiple AI agents, where value is generated through complex interactions and emergent behaviors, moving beyond simple per-task or per-token billing.

These research directions collectively point toward a richer, more nuanced understanding of AI pricing and its implications for theory, practice, and policy. Addressing these questions will be crucial for ensuring the sustainable, equitable, and responsible development of the AI-driven economy.

---

## Conclusion

The advent of artificial intelligence (AI) has ushered in a transformative era for business and economic landscapes, fundamentally reshaping how firms conceptualize, implement, and manage pricing strategies (Liu, 2024)(Obiajulu et al., 2025). This paper embarked on a comprehensive theoretical analysis, complemented by insights from various case studies, to explore the intricate dynamics of AI pricing, its adoption patterns, and the emergent potential for collusive behaviors in digitally mediated markets. Our investigation has illuminated how AI-driven pricing mechanisms, while offering unprecedented efficiency and optimization capabilities, also introduce novel complexities and regulatory challenges that demand careful consideration from both academic and practical standpoints (Mirghaderi et al., 2023)(Ganguly, 2025). The core objective was to dissect the multifaceted implications of AI integration into pricing models, particularly focusing on its capacity to foster more sophisticated, dynamic, and potentially opaque market interactions, thereby influencing competitive structures and consumer welfare.

The theoretical framework applied throughout this study underscored the shift from traditional, human-centric pricing decisions to algorithms that learn, adapt, and execute

pricing adjustments in real-time, often with minimal human intervention (Niharika et al., 2024)(Vashishtha et al., 2022). We demonstrated that AI's ability to process vast quantities of data–including competitor pricing, consumer behavior, inventory levels, and external market conditions–allows for hyper-segmentation and personalized pricing strategies that were previously unattainable (Shiva Kumar Bhuram, 2025)(Subham, 2025). This dynamic pricing capability, while maximizing revenue and optimizing resource allocation for individual firms, simultaneously creates a complex interplay of interdependent pricing decisions across an industry (Ye & Zhang, 2017)(Sonderegger, 2011). The analysis revealed that the adoption of AI pricing is not merely an incremental technological upgrade but a strategic transformation that can fundamentally alter market equilibrium and competitive intensity. Firms leveraging AI gain a significant competitive advantage by optimizing price elasticity of demand (Brynjolfsson et al., 2023), managing inventory more effectively (Abbas, 2025), and responding to market shifts with unparalleled agility (Divakaruni & Navarro, 2024). The economic benefits for adopting firms, ranging from enhanced profitability to improved market share, often drive a rapid diffusion of these technologies, creating a "race to the bottom" or a "race to the top" depending on market concentration and competitive dynamics.

A central finding of this research pertains to the heightened risk of algorithmic collusion, an emergent phenomenon where AI-driven pricing algorithms, operating independently, converge on tacit collusive outcomes without explicit human communication or agreement (Liu, 2024). Our theoretical models illustrated how algorithms, designed to maximize individual firm profits, can learn to anticipate and respond to competitor pricing actions in a way that leads to supra-competitive prices, mimicking the effects of traditional cartels. This self-organizing collusion poses significant challenges for antitrust enforcement, as the lack of direct communication makes it difficult to prove intent or coordinated action (Lo, 2018). The study highlighted various mechanisms through which such collusion could manifest, including parallel adoption of similar pricing algorithms (Diaw & Pouyet, 2004), learning through repeated interactions in oligopolistic markets, and the use of common third-party

pricing software that, inadvertently or by design, facilitates price coordination (Cody, 2000). The analysis suggests that the transparency and interpretability of these algorithms become paramount in mitigating such risks, yet many advanced AI systems, particularly large language models (LLMs) used in pricing analytics (Barbere et al., 2024)(Rudnytskyi, 2022), are often characterized by their "black box" nature, making it difficult to discern their internal decision-making processes (Mirghaderi et al., 2023). This opacity further complicates regulatory oversight and the development of effective detection mechanisms.

This paper offers several significant contributions to both theoretical understanding and practical application in the fields of business, economics, and information systems. Theoretically, we extend the literature on industrial organization by integrating the unique characteristics of AI into models of firm behavior and market competition. We provide a nuanced framework for understanding how AI's capabilities–such as real-time data processing, predictive analytics (Niharika et al., 2024)(Subham, 2025), and autonomous decision-making– interact with existing market structures to produce novel competitive outcomes, including the potential for algorithmic collusion. This moves beyond traditional analyses of explicit collusion to address the more subtle and emergent forms of coordination facilitated by advanced algorithms. Furthermore, by exploring the value creation and capture mechanisms in AI-driven ecosystems (Lorente, 2025), we contribute to the strategic management literature, illustrating how firms can derive and sustain competitive advantage through the strategic deployment of AI in pricing, while also acknowledging the associated risks. The research also sheds light on the evolving nature of transaction cost economics (Williamson, 2010) in digital markets, where AI can significantly reduce information asymmetries and search costs, but simultaneously introduce new forms of market power and lock-in effects (Vogt & Bizer, 2013).

From a practical perspective, the findings offer critical insights for managers, policy-makers, and regulators. For managers, the study underscores the strategic imperative of AI adoption in pricing but also highlights the need for robust governance frameworks around algorithmic transparency and ethical considerations (Ganguly, 2025). Understanding the

potential for algorithmic collusion is crucial for firms to navigate competitive landscapes responsibly and avoid inadvertent anti-competitive practices. The research emphasizes that while AI offers immense opportunities for revenue optimization and efficiency (Kshirsagar et al., 2021)(Cho & Bahn, 2020), its implementation requires a deep understanding of market dynamics and a commitment to fair competition. For policymakers and antitrust regulators, this paper serves as a clarion call for developing new regulatory tools and frameworks capable of addressing the unique challenges posed by AI pricing (Lo, 2018). Traditional antitrust laws, designed for human-led collusion, may prove inadequate in detecting and prosecuting algorithmic collusion. New approaches, potentially focusing on algorithm design, data inputs, and real-time market monitoring, are necessary to safeguard consumer welfare and maintain competitive markets (Mirghaderi et al., 2023). The discussion on API monetization (De, 2017)(Ladas et al., 2019) and cloud service pricing (Satapathi, 2025)(Livingstone, 2013) further provides actionable intelligence for businesses operating in platform economies, emphasizing the need for transparent and fair pricing models that reflect the true value delivered (Obiajulu et al., 2025).

Despite these contributions, this study is subject to certain limitations that offer fertile ground for future research. Primarily, as a theoretical analysis, the generalizability of some findings may be constrained by the assumptions embedded in our models. While case studies provided illustrative examples, comprehensive empirical validation across a diverse range of industries and market structures is still needed to fully substantiate the theoretical propositions. The "black box" nature of many advanced AI algorithms also presents a challenge, as our analysis largely relies on inferring algorithmic behavior rather than direct observation of their internal decision-making processes. Moreover, the rapid evolution of AI technology means that the specific mechanisms of algorithmic collusion and their detection methods are continually changing, requiring ongoing research and adaptation.

Future research should therefore prioritize empirical studies that leverage real-world data to test the propositions advanced in this paper. This could involve analyzing pricing data

from industries with high AI adoption rates, conducting controlled experiments with AI pricing agents, or developing more sophisticated econometric models to detect patterns indicative of algorithmic collusion (Liu, 2024). Further exploration into the design of "collusion-resistant" AI algorithms or regulatory sandboxes for testing AI pricing strategies could also yield valuable insights. Additionally, interdisciplinary research combining economic theory with computer science and legal scholarship is crucial for developing robust regulatory frameworks and technical solutions to ensure ethical and fair AI pricing practices (Ganguly, 2025). The role of human oversight in AI pricing, the development of explainable AI (XAI) for pricing algorithms, and the implications of AI for consumer privacy and fairness in personalized pricing warrant deeper investigation (Fang & Zhou, 2025)(Siddannavar et al., 2025)(Yin & Qiu, 2021). As AI continues to permeate various sectors, from automotive aftermarkets (Shiva Kumar Bhuram, 2025) to financial sentiment analysis (Leechewyuwasorn & Wangpratham, 2024) and resource allocation (Li & Ren, 2025)(Huang et al., 2024), understanding its pervasive impact on pricing and market competition will remain a critical area of academic inquiry and practical concern. The journey towards harnessing AI's full potential while mitigating its risks is ongoing, and this paper serves as a foundational step in that essential endeavor.

---

# Appendix A: Framework for AI Pricing Model Design and Evaluation

*A.1 Introduction to the Framework*

This appendix details a comprehensive framework for designing, evaluating, and optimizing pricing models for Artificial Intelligence (AI) services. Building upon the multi-dimensional theoretical framework presented in the Methodology section, this expanded framework provides a structured approach for practitioners and researchers to systematically

assess how different pricing strategies align with business objectives, technological constraints, market dynamics, and ethical considerations. The increasing complexity and diversity of AI applications necessitate a nuanced approach that transcends simple cost-plus or competitive pricing, moving towards models that effectively capture the dynamic value generated by AI while ensuring sustainability and fairness.

The framework is designed to be iterative and adaptive, recognizing the rapid evolution of AI technologies and market conditions. It integrates elements of economic theory, strategic management, and ethical AI principles, offering a holistic perspective on AI monetization. The core objective is to guide stakeholders in developing pricing models that are not only economically viable but also promote widespread adoption, foster trust, and contribute to the responsible development of the AI ecosystem.

*A.2 Core Components of the Framework*

The framework is structured around five interconnected phases, each comprising specific considerations and analytical tools. These phases guide the user from an initial understanding of the AI service to its ongoing optimization in the market.

**A.2.1 Phase 1: AI Service Definition and Value Proposition Articulation** This initial phase focuses on clearly defining the AI service, its capabilities, and the unique value it delivers to target customers. A precise understanding of the value proposition is foundational for any effective pricing strategy.

- **A.2.1.1 Service Capabilities and Features:**
- What specific tasks does the AI perform (e.g., text generation, image recognition, predictive analytics)?
- What are its key features, performance metrics (accuracy, latency, throughput), and unique selling points?
- Is it a foundational model, a specialized agent, or an embedded AI component?

- **A.2.1.2 Target Customer Segments:**

- Who are the primary users (e.g., individual developers, startups, SMEs, large enterprises, academic researchers)?

- What are their specific needs, pain points, and existing workflows?

- How do their technical capabilities, budget constraints, and willingness to pay vary?

- **A.2.1.3 Quantifiable Value Generation:**

- How does the AI service create tangible economic benefits (e.g., cost savings, revenue increase, efficiency gains, risk reduction)?

- What are the key performance indicators (KPIs) that customers will use to measure the AI's impact?

- Can the value be directly attributed to the AI, or is it an indirect contribution to a broader process?

- **A.2.1.4 Strategic Differentiation:**

- How does the AI service stand out from competitors or alternative solutions (e.g., superior performance, ease of integration, ethical design, specialized domain knowledge)?

- What are the competitive advantages that justify a specific price point?

**A.2.2 Phase 2: Cost Structure Analysis and Sustainability Assessment** This phase involves a detailed analysis of all costs associated with developing, deploying, and maintaining the AI service, ensuring the pricing model supports long-term sustainability.

- **A.2.2.1 Fixed Costs:**

- Research & Development (R&D) investments (model training, data acquisition, human capital).

- Infrastructure setup (GPU clusters, data centers).

- Software licensing and tooling.

- **A.2.2.2 Variable Costs:**

- Computational resources (GPU/CPU hours for inference, memory, storage).

- Data transfer (ingress/egress).

- API calls to third-party services.

- Ongoing data labeling and model fine-tuning.

- Customer support and operational overhead per user.

- **A.2.2.3 Cost Attribution and Recovery:**

- How can variable costs be accurately attributed to individual user interactions or service units (e.g., per token, per request)?

- What is the desired profit margin, and how will it be achieved across different service tiers?

- Are there economies of scale or scope that can reduce marginal costs over time?

- **A.2.2.4 Green AI Costs:**

- What is the energy consumption and carbon footprint of the AI model's training and inference?

- How can these environmental costs be internalized or reflected in pricing to incentivize sustainable usage?

**A.2.3 Phase 3: Market & Competitive Landscape Analysis**   This phase examines the external market environment, including competitive offerings, demand characteristics, and regulatory factors.

- **A.2.3.1 Competitive Analysis:**

- Who are the direct and indirect competitors?

- What are their pricing models, feature sets, and market positioning?

- What are their strengths and weaknesses?

- **A.2.3.2 Price Elasticity of Demand:**

- How sensitive are different customer segments to price changes for this AI service?

- Are there substitutes or alternative solutions that influence demand elasticity?

- **A.2.3.3 Market Segmentation and Positioning:**

- How can the market be segmented based on needs, usage patterns, and willingness to pay?

- What is the desired market positioning (e.g., premium, value, mass market)?

- **A.2.3.4 Regulatory Environment:**

- What existing or emerging regulations (e.g., data privacy, antitrust, AI governance) might impact pricing flexibility or require specific transparency measures?

- Are there concerns about market power or algorithmic collusion?

**A.2.4 Phase 4: Pricing Model Selection and Design**  Based on the insights from the previous phases, this phase involves selecting and designing the most appropriate pricing model(s).

- **A.2.4.1 Model Selection:**

- Consider pure models (token-based, request-based, subscription, compute-based, value-based).

- Evaluate the suitability of hybrid models (e.g., subscription + usage, feature-gated + usage).

- Explore dynamic pricing mechanisms (real-time adjustments based on demand, supply, user segment).

- **A.2.4.2 Tiered Structure Design:**

- Define different service tiers (e.g., Free, Basic, Pro, Enterprise).

- Specify features, usage limits, performance guarantees (SLAs), and support levels for each tier.

- Determine the pricing points and value differentials between tiers.

- **A.2.4.3 Billing Metrics and Granularity:**

- Choose appropriate billing units (e.g., per token, per API call, per GPU hour, per outcome).

- Determine the level of granularity required for fairness and cost alignment.

- Decide on differential pricing for input/output, model versions, or specific features.

- **A.2.4.4 Contractual Terms and Policies:**

- Define terms of service, fair-use policies, and overage charges.

- Establish clear mechanisms for value attribution in outcome-based models.

- Address data usage policies and privacy guarantees.

**A.2.5 Phase 5: Implementation, Monitoring, and Iteration** The final phase focuses on the practical aspects of deploying the pricing model, continuously monitoring its performance, and iterating based on feedback and market changes.

- **A.2.5.1 Billing System Integration:**

- Ensure robust and accurate billing infrastructure that can track usage across all defined metrics.

- Provide clear and transparent billing statements to customers.

- **A.2.5.2 Performance Monitoring:**

- Track key metrics: revenue per user, customer acquisition cost, churn rate, usage patterns, resource consumption.

- Monitor competitor pricing and market trends.

- Gather customer feedback on pricing fairness and transparency.

- **A.2.5.3 Iteration and Optimization:**

- Regularly review the pricing model against defined objectives.

- Be prepared to adjust pricing points, tiers, or models in response to market changes, technological advancements, or customer feedback.

- Leverage AI-driven analytics to continuously optimize pricing strategies.

- **A.2.5.4 Ethical and Regulatory Compliance:**

- Continuously ensure the pricing model adheres to all relevant ethical guidelines and legal regulations.

- Proactively address any concerns regarding fairness, transparency, or potential anti-competitive effects.

*A.3 Conclusion of the Framework*

This comprehensive framework provides a structured, iterative approach to navigating the complexities of AI pricing. By systematically addressing AI service definition, cost structures, market dynamics, pricing model design, and ongoing optimization, organizations can develop pricing strategies that are not only economically sustainable but also foster widespread adoption, build customer trust, and contribute to the responsible growth of the AI industry. The framework emphasizes that AI pricing is a continuous strategic process, requiring adaptability and a deep understanding of both the technology and its market context.

---

# Appendix C: Detailed Case Study Projections and Quantitative Metrics

This appendix provides detailed quantitative metrics and illustrative projections for various scenarios related to AI pricing and adoption, drawing from the real-world implementations discussed in the Analysis section. These tables aim to further elucidate the economic implications of different pricing models and the value generated by AI services, providing a more granular view for strategic planning and cost optimization.

*C.1 Scenario 1: OpenAI GPT-4 API Usage Cost Projections*

This scenario projects the monthly costs for an enterprise utilizing OpenAI's GPT-4 API for content generation, customer support automation, and code assistance, comparing costs across different usage volumes. It highlights the impact of token-based pricing on variable workloads.

## Table C.1: OpenAI GPT-4 API Monthly Cost Projections (Illustrative)

| Use Case | Avg. Input Tokens/Req | Avg. Output Tokens/Req | Req/Month (Low) | Cost/Month (Low) | Req/Month (High) | Cost/Month (High) |
|---|---|---|---|---|---|---|
| **Content Gen.** | 500 | 1500 | 10,000 | $225.00 | 100,000 | $2,250.00 |
| **Customer Support** | 200 | 300 | 50,000 | $1,125.00 | 500,000 | $11,250.00 |
| **Code Assist** | 1000 | 800 | 5,000 | $112.50 | 50,000 | $1,125.00 |
| **Total Est. Cost** | | | | **$1,462.50** | | **$14,625.00** |

*Note: Assumes GPT-4-8k pricing: $0.03/1k input tokens, $0.06/1k output tokens. Costs are illustrative and subject to change by provider. "Req" = Requests.*

*C.2 Scenario 2: Anthropic Claude 3 Opus vs. Sonnet Cost-Efficiency*

This scenario compares the cost-efficiency of Anthropic's Claude 3 Opus (premium) and Claude 3 Sonnet (mid-tier) models for processing large documents, focusing on how context window and token pricing impact total cost for a fixed task (e.g., summarizing a 50,000-token document).

## Table C.2: Claude 3 Model Cost Comparison for Document Summarization

| Model | Input Tokens/1k (USD) | Output Tokens/1k (USD) | Document Size (Tokens) | Summary Output (Tokens) | Total Input Cost | Total Output Cost | Total Cost |
|---|---|---|---|---|---|---|---|
| **Claude 3 Opus** | $15.00 | $75.00 | 50,000 | 5,000 | $0.75 | $0.375 | $1.125 |

| Model | Input Tokens/1k (USD) | Output Tokens/1k (USD) | Document Size (Tokens) | Summary Output (Tokens) | Total Input Cost | Total Output Cost | Total Cost |
|---|---|---|---|---|---|---|---|
| **Claude 3 Sonnet** | $3.00 | $15.00 | 50,000 | 5,000 | $0.15 | $0.075 | $0.225 |

*Note: Pricing is illustrative based on published rates. Assumes a 10% summary output. Opus is 5x more expensive per token than Sonnet, but may offer higher quality for complex tasks. This table highlights the trade-off between cost and potential quality/capability.*

C.3 Scenario 3: Enterprise Cloud AI Cost Optimization via Hybrid Models

This scenario illustrates how an enterprise might optimize its AI spending for a custom LLM deployment using a hybrid cloud computing and managed AI service model. It projects costs for training a model on dedicated compute instances versus using a managed fine-tuning service.

**Table C.3: Custom LLM Deployment Cost Optimization (Illustrative Annual Costs)**

| Cost Category | Dedicated Compute (USD) | Managed Fine-tuning (USD) | Hybrid Approach (USD) | Optimization Impact |
|---|---|---|---|---|
| **GPU Instances (Training)** | $250,000 | $180,000 | $180,000 | 28% cost reduction |
| **GPU Instances (Inference)** | $150,000 | $150,000 | $100,000 | 33% cost reduction |

| Cost Category | Dedicated Compute (USD) | Managed Fine-tuning (USD) | Hybrid Approach (USD) | Optimization Impact |
|---|---|---|---|---|
| **Data Storage & Transfer** | $30,000 | $25,000 | $20,000 | 33% cost reduction |
| **Managed Service Fees** | $0 | $50,000 | $30,000 | Variable |
| **Developer/Ops Overhead** | $100,000 | $50,000 | $40,000 | 60% reduction |
| **Total Annual Cost** | **$530,000** | **$455,000** | **$370,000** | **30% overall savings** |

*Note: "Dedicated Compute" assumes full management by the enterprise. "Managed Fine-tuning" uses a provider's service. "Hybrid Approach" combines managed services for training with optimized, serverless inference. Illustrative figures only.*

*C.4 Scenario 4: Value-Based Pricing for AI-Driven Fraud Detection*

This scenario projects the potential financial impact of adopting an AI-driven fraud detection system priced on a value-based model (e.g., a percentage of detected fraud savings). It shows how a provider's revenue is tied to the value delivered.

**Table C.4: AI-Driven Fraud Detection: Value-Based Pricing Projections**

| Metric | Baseline (No AI) | AI System (Tier 1) | AI System (Tier 2) | AI System (Tier 3) |
|---|---|---|---|---|
| **Annual Fraud Losses** | $10,000,000 | $5,000,000 | $3,000,000 | $1,500,000 |
| **Fraud Loss Reduction (%)** | 0% | 50% | 70% | 85% |

| Metric | Baseline (No AI) | AI System (Tier 1) | AI System (Tier 2) | AI System (Tier 3) |
|---|---|---|---|---|
| **Annual Savings (USD)** | $0 | $5,000,000 | $7,000,000 | $8,500,000 |
| **AI Provider Fee (%)** | N/A | 10% | 8% | 6% |
| **AI Provider Revenue (USD)** | $0 | $500,000 | $560,000 | $510,000 |
| **Net Client Savings (USD)** | $0 | $4,500,000 | $6,440,000 | $7,990,000 |

*Note: This projection illustrates a value-based pricing model where the AI provider's fee is a percentage of the fraud losses prevented. Different tiers represent varying levels of AI sophistication and corresponding fraud reduction capabilities, with potentially decreasing percentage fees for higher savings to incentivize adoption.*

## C.5 Conclusion of Detailed Projections

These detailed projections underscore the complex interplay of pricing models, usage patterns, and value realization in the AI ecosystem. They demonstrate that while token-based and usage-based models offer granular control and cost alignment, value-based and hybrid approaches hold significant potential for optimizing both provider revenue and customer ROI. The ability to accurately forecast costs, manage resource consumption, and clearly articulate the economic benefits of AI solutions will be paramount for successful adoption and sustainable growth in this rapidly evolving market. These quantitative perspectives provide a crucial foundation for both theoretical understanding and practical decision-making for stakeholders navigating the AI pricing landscape.

# Appendix D: Additional References and Resources

This appendix provides a curated list of supplementary resources that delve deeper into the foundational theories, technological advancements, and practical implications discussed in the main thesis. These resources are categorized to facilitate targeted exploration for readers seeking further information on specific aspects of AI pricing and the broader AI ecosystem.

*D.1 Foundational Texts in Digital Economics and Pricing*

1. **Shapiro, C., & Varian, H. R. (1999).** *Information Rules: A Strategic Guide to the Network Economy.* **Harvard Business Review Press.**

- **Relevance:** This seminal work provides foundational economic principles for understanding pricing, competition, and strategy in the information age, highly relevant to digital goods like AI services.

2. **Varian, H. R. (1995).** *Microeconomic Analysis* **(3rd ed.). W. W. Norton & Company.**

- **Relevance:** A classic microeconomics textbook offering in-depth coverage of pricing theory, market structures, and consumer behavior, essential for understanding the underlying principles of AI pricing.

3. **Nagle, T. T., & Müller, G. (2017).** *The Strategy and Tactics of Pricing: A Guide to Growing More Profitably* **(6th ed.). Pearson Education.**

- **Relevance:** Provides practical insights into various pricing strategies, including value-based pricing, and how to implement them effectively in a competitive market.

*D.2 Key Research Papers on AI/LLM Economics*

1. **Agrawal, A., Gans, J. S., & Goldfarb, A. (2018).** *Prediction Machines: The Simple Economics of Artificial Intelligence.* **Harvard Business Review Press.**

- **Relevance:** Explores the economic implications of AI as a "prediction technology," offering insights into how businesses can leverage and price AI capabilities.

2. **Korinek, A., & Stiglitz, J. E. (2021).** *Artificial Intelligence and Economic Growth: Automation, Distribution, and the Need for Policy Adjustment.* **National Bureau of Economic Research (NBER) Working Paper 28625.**

- **Relevance:** Discusses the broader economic impact of AI, including its effects on productivity, labor markets, and the distribution of wealth, providing context for AI monetization strategies.

3. **Etzioni, A., & Etzioni, O. (2017).** *AI Assisted Ethics.* **AI & Society, 32(2), 273-278.**

- **Relevance:** While not directly on pricing, this paper on AI ethics is crucial for understanding the ethical considerations that increasingly influence AI governance and pricing transparency.

*D.3 Online Resources and Industry Reports*

- **OpenAI Pricing Page**: https://openai.com/api/pricing - Direct source for OpenAI's token-based pricing for their various LLMs.

- **Anthropic Pricing Page**: https://www.anthropic.com/api/pricing - Provides detailed token pricing for Anthropic's Claude models, often highlighting large context windows.

- **Google Cloud AI Pricing**: https://cloud.google.com/vertex-ai/pricing - Details pricing for Google's AI platform services, including generative AI models.

- **Microsoft Azure AI Pricing**: https://azure.microsoft.com/en-us/pricing/details/cognitive-services/ - Provides pricing for Azure's comprehensive suite of AI services.

- **AWS Bedrock Pricing**: https://aws.amazon.com/bedrock/pricing/ - Outlines consumption-based pricing for foundational models offered through AWS Bedrock.

- **Hugging Face Inference Endpoints**: https://huggingface.co/docs/inference-endpoints/pricing - Details compute-based pricing for deploying open-source models as managed services.

- **Gartner / Forrester / IDC Reports (various)**: Access through institutional subscriptions. These reports offer market analyses, competitive landscapes, and future predictions for AI and cloud services, often including pricing trends.

*D.4 Software and Tools for AI Cost Management*

- **Cloud Cost Management Platforms (e.g., FinOps tools)**: Solutions like Cloud-Health by VMware, Flexera, or Apptio provide capabilities for monitoring, optimizing, and forecasting cloud and AI expenditures.

- **LLM Tokenizers**: Online tools and libraries (e.g., OpenAI's tiktoken, Hugging Face tokenizers) allow users to estimate token counts for text, aiding in cost prediction for token-based pricing.

- **API Management Platforms**: Tools like Apigee (Google), Azure API Management, or AWS API Gateway assist in monitoring API usage, applying rate limits, and potentially managing tiered access, which can inform usage-based billing.

*D.5 Professional Organizations and Communities*

- **AI Ethics & Governance Forum**: Various academic and industry groups focusing on responsible AI development, which often touch upon fair pricing and transparency.

- **FinOps Foundation**: A community that advances the practice of cloud financial management, highly relevant for optimizing the costs of AI services deployed in cloud environments.

- **Open-Source AI Communities (e.g., Hugging Face, GitHub)**: Platforms where open-source LLMs are developed and discussed, influencing the competitive landscape and de facto pricing floors for commercial models.

# Appendix E: Glossary of Terms

This glossary defines key technical terms and domain-specific jargon used throughout the thesis, providing clear and concise explanations to enhance reader comprehension.

**Agentic AI**: An artificial intelligence system designed to exhibit autonomous, goal-oriented behavior, capable of planning, reasoning, and executing complex tasks across various environments.

**Algorithmic Collusion**: An emergent phenomenon where independent AI pricing algorithms, designed to maximize individual firm profits, converge on tacit collusive outcomes without explicit human communication or agreement, leading to supra-competitive prices.

**API Call**: A request made to an Application Programming Interface (API) to access a specific service or functionality; often a unit of billing in usage-based pricing models for AI services.

**Attribution Problem**: The challenge of precisely determining how much of a specific business outcome (e.g., revenue increase, cost savings) can be directly and solely attributed to an AI solution, rather than other influencing factors.

**Black Box AI**: Refers to AI systems, particularly complex deep learning models, whose internal decision-making processes are opaque and difficult for humans to understand or interpret.

**Cloud Computing**: The delivery of on-demand computing services–including servers, storage, databases, networking, software, analytics, and intelligence–over the Internet ("the cloud") with pay-as-you-go pricing.

**Compute-Based Pricing**: A pricing model where users are charged directly for the underlying computational resources consumed (e.g., CPU/GPU hours, memory, storage) by their AI workloads.

**Context Window**: The maximum number of tokens (input + output) that a Large Language Model can process or "pay attention to" at one time during an interaction.

**Cost-Plus Pricing**: A traditional pricing strategy where a markup is added to the total cost of producing a product or service to determine its selling price.

**Customer Lifetime Value (CLV)**: A prediction of the total revenue a business can expect to generate from a customer throughout their relationship.

**Dynamic Pricing**: A pricing strategy where prices are adjusted in real-time based on fluctuating factors such as demand, supply, time of day, customer segment, and competitor actions.

**Edge Computing**: A distributed computing paradigm that brings computation and data storage closer to the sources of data, often to reduce latency and bandwidth usage.

**Fair-Use Policy**: A set of rules or guidelines specifying acceptable usage limits for a service, often employed in subscription models to prevent excessive consumption by a small number of users.

**Feature-Based Pricing**: A pricing model that differentiates costs based on the specific capabilities, performance levels, or advanced functionalities offered by a product or service.

**Foundational Model**: A large AI model, typically a Large Language Model, trained on a vast amount of data, capable of being adapted to a wide range of downstream tasks.

**Freemium Model**: A business strategy where a basic version of a product or service is offered for free to attract users, with premium features or higher usage limits available through a paid subscription.

**Generative AI**: A type of artificial intelligence that can produce various types of content, including text, images, audio, and synthetic data.

**Green AI**: An approach to AI development and deployment that prioritizes energy efficiency and environmental sustainability, aiming to reduce the carbon footprint of AI operations.

**Hybrid Pricing Model**: A pricing strategy that combines elements from two or more basic pricing paradigms (e.g., subscription with usage-based overages) to optimize for diverse user needs and business objectives.

**Inference**: The process of using a trained AI model to make predictions or generate outputs based on new input data.

**Large Language Model (LLM)**: A type of artificial intelligence algorithm that uses deep learning techniques and massive datasets to understand, summarize, generate, and predict new content.

**Marginal Cost**: The cost incurred by producing one additional unit of a good or service; for digital goods, this often approaches zero after initial development.

**Multimodal AI**: AI systems capable of processing and generating information across multiple modalities, such as text, images, audio, and video.

**Non-Rivalrous Good**: A good whose consumption by one person does not prevent others from consuming it simultaneously; a characteristic of many digital goods.

**Outcome-Based Pricing**: A pricing model where payment is directly tied to the achievement of specific, measurable business results or performance metrics delivered by the AI solution.

**Pay-Per-Use**: A pricing model where customers pay only for the actual amount of a service or resource they consume, without fixed recurring fees.

**Price Elasticity of Demand**: A measure of the responsiveness of the quantity demanded of a good or service to a change in its price.

**Prompt Engineering**: The process of carefully crafting input queries (prompts) for generative AI models to elicit desired and optimal responses, often to manage token usage and cost.

**Retrieval-Augmented Generation (RAG)**: An AI architecture that enhances LLM performance by retrieving relevant information from an external knowledge base and

incorporating it into the prompt, rather than relying solely on the model's pre-trained knowledge.

**Subscription-Based Pricing**: A pricing model where customers pay a recurring fixed fee (e.g., monthly or annually) for access to a service, often with predefined usage limits or feature sets.

**Token**: A discrete unit of text (e.g., a word, sub-word, punctuation mark) into which a Large Language Model breaks down input and output text for processing.

**Token-Based Pricing**: A granular usage-based pricing model specifically tailored to LLMs, where users are charged based on the number of tokens consumed for both input prompts and generated output.

**Transaction Cost Economics (TCE)**: A framework that analyzes the costs associated with making economic exchanges, including search, bargaining, and enforcement costs, influencing "make or buy" decisions.

**Transparency (in Pricing)**: The clarity and comprehensibility of a pricing structure, allowing users to easily understand how costs are calculated and predict their expenditures.

**Usage-Based Pricing**: A broad category of pricing models where customers are charged based on their consumption of various resources or features (e.g., API calls, compute hours, data processed).

**Value-Based Pricing**: A pricing strategy that sets prices primarily based on the perceived or actual economic value delivered to the customer, rather than on production costs or competitor prices.

**Versioning**: A pricing strategy that involves offering different qualities, feature sets, or performance levels of a product or service at various price points to cater to diverse customer segments.

---

# References

Abbas. (2025). Optimizing Surplus Inventory Management in Electrical Distribution Through API-Driven Marketplace Integration. *Entrepreneurship and Management Science Journal.* https://doi.org/10.59573/emsj.9(4).2025.29.

Anonymous. (2025). AI-Driven Strategies for Cloud Cost Optimization. *International Journal of Scientific & Advanced Technology.* https://doi.org/10.71097/ijsat.v16.i2.4714.

Ayata. (2020). Old abuses in new markets? Dealing with excessive pricing by a two-sided platform. *Journal of Antitrust Enforcement.* https://doi.org/10.1093/jaenfo/jnaa008.

Barbere, Martin, Thornton, Harris, & Thompson. (2024). *Dynamic Token Hierarchies: Enhancing Large Language Models with a Multi-Tiered Token Processing Framework.* https://doi.org/10.36227/techrxiv.172971998.83622138/v1

Brynjolfsson, Jin, & Unger. (2023). *The Price Elasticity of Demand for AI Software: Evidence from the Cloud.* National Bureau of Economic Research (NBER). https://doi.org/10.3386/w31751

Chaerul. (2025). Analyzing User Expectation and Experience to Formulate Data-Driven Strategy for Elevating Paid Subscription Engagement in Ai Educational Tools. *Enrichment: Journal of Management.* https://doi.org/10.55324/enrichment.v3i4.452.

Chiruvelli. (2025). Generative AI For Personalized Product Bundling In Consumer Banking: A Revenue Optimization Framework. *Journal of Innovative Computing, Research, and Reviews.* https://doi.org/10.63278/jicrcr.vi.3302.

Cho, & Bahn. (2020). A Cost Estimation Model for Cloud Services and Applying to PC Laboratory Platforms. *Processes.* https://doi.org/10.3390/pr8010076.

Cody. (2000). *Competitive Infrastructure: As An Enabler of Market-Based Pricing.* Springer. https://doi.org/10.1007/978-1-4615-4529-3_4

De. (2017). *API Monetization.* Springer. https://doi.org/10.1007/978-1-4842-1305-6_8

Diaw, & Pouyet. (2004). *Competition, Incomplete Discrimination and Versioning.* https://doi.org/10.2139/ssrn.606961

Divakaruni, & Navarro. (2024). *Technology Adoption and Pricing: Evidence from US Airlines.* https://doi.org/10.2139/ssrn.4718902

Fang, & Zhou. (2025). *Understanding the Impacts of Human-Like Competencies on Users' Willingness to Pay for Ai Companion Services: A Mixed-Method Approach.* https://doi.org/10.2139/ssrn.5333712

Fishburn, & Odlyzko. (1999). Competitive pricing of information goods: Subscription pricing versus pay-per-use. *Economic Theory.* https://doi.org/10.1007/s001990050264.

Ganguly. (2025). *AI Governance and Responsible AI.* Springer. https://doi.org/10.1007/979-8-8688-1154-8_7

Gbadamosi, Nwulu, & Sun. (2018). Harmonic and power loss minimization in power systems incorporating renewable energy sources and locational marginal pricing. https://doi.org/10.1063/1.5041923

Geetha, Ayub, Vivin, & Chandran. (2024). THE INFLUENCE OF ADOPTING ARTIFICIAL INTELLIGENCE (AI) ON MALAYSIA'S ECONOMIC ENVIRONMENT.. *Malaysian Journal of Business and Economics.* https://doi.org/10.51200/mjbe.v11i1.5294.

Huang, Li, Yang, Si, Ma, & Wang. (2024). Multiparticipant Double Auction for Resource Allocation and Pricing in Edge Computing. *IEEE Internet of Things Journal.* https://doi.org/10.1109/JIOT.2023.3339655.

Kaaniche, & Laurent. (2018). BDUA: Blockchain-Based Data Usage Auditing. https://doi.org/10.1109/cloud.2018.00087

Korinek. (2025). *AI Agents for Economic Research.* National Bureau of Economic Research (NBER). https://doi.org/10.3386/w34202

Kshirsagar, More, Lahoti, Adgaonkar, Jain, & Ryan. (2021). GREE-COCO: Green Artificial Intelligence Powered Cost Pricing Models for Congestion Control. https://doi.org/10.5220/0010261209160923

Kärrberg. (2010). A Two-Sided Market Approach to Value Chain Dynamics in Telecom Services: A Study Lens for Mobile Platform Innovation and Pricing Strategies. https://doi.org/10.1109/ICMB-GMR.2010.76

Ladas, Kavadias, & Loch. (2019). *Product Selling Versus Pay-Per-Use Services: A Strategic Analysis of Competing Business Models.* https://doi.org/10.2139/ssrn.3356458

Leechewyuwasorn, & Wangpratham. (2024). *Comparative Analysis of Financial Sentiment Analysis Models for the Thai Stock Market: Traditional NLP vs. GPT vs. Gemini.* https://doi.org/10.2139/ssrn.4921837

Li, & Ren. (2025). Enhancing Multi-Tenant Database Resource Allocation with Active Learning-Driven Multi-Objective Optimization. https://doi.org/10.1109/eeiss65394.2025.11085993

Liu. (2024). *AI Pricing: Adoption of Artificial Intelligences and Collusive Price.* https://doi.org/10.62051/ed63gf48

Livingstone. (2013). *Cloud price wars could drive 'volatility as a service'.* https://doi.org/10.64628/aa.fwtdh7kfg

Lo. (2018). MARKETER'S PRICING STRATEGY VS. COMPETITION LAW. *Lincoln Journal of Management and Social Sciences.* https://doi.org/10.51200/ljms.v12i.1353.

Lorente. (2025). Value Creation and Value Capture in AI: A Triple Helix Model. *AI Ethics.* https://doi.org/10.1609/aies.v8i2.36662.

Ma, Yang, Huang, Li, & Wang. (2015). The Re-Engineering of Iron and Steel Enterprise Procurement Processes Based on the ABC Method. ASCE. https://doi.org/10.1061/9780784479384.115

Maguire. (2021). *Value selling.* Routledge. https://doi.org/10.4324/9781003177937-20

Mirghaderi, Sziron, & Hildt. (2023). Ethics and Transparency Issues in Digital Platforms: An Overview. *AI.* https://doi.org/10.3390/ai4040042.

Nagaraju, Bahrami, Prasad, Mantena, Biswal, & Islam. (2023). Predicting California Bearing Ratio of Lateritic Soils Using Hybrid Machine Learning Technique. *Buildings.* https://doi.org/10.3390/buildings13010255.

Niharika, Hareesh, & Ariwa. (2024). *Pricing Optimisation Using Predictive Analytics.* CRC Press. https://doi.org/10.1201/9781003472544-8

Obiajulu, Azubuko, & Esemokhai. (2025). Pricing transparency and innovative pricing models: Driving customer trust and retention in SMEs. *World Journal of Advanced Research and Reviews.* https://doi.org/10.30574/wjarr.2025.27.3.2924.

Oleksii. (2025). Algorithmizing B2B Sales: Can AI Create a Sales Framework That Guarantees Predictable Results?. *The American Journal of Management and Economics Innovations.* https://doi.org/10.37547/tajmei/volume07issue03-02.

Rodeghero, McMillan, & Shirey. (2017). API Usage in Descriptions of Source Code Functionality. https://doi.org/10.1109/wapi.2017.3

Rudnytskyi. (2022). *openai: R Wrapper for OpenAI API.* https://doi.org/10.32614/cran.package.openai

Satapathi. (2025). *Pricing tiers of Azure AI Language Service.* Springer. https://doi.org/10.1007/979-8-8688-1333-7_4

Seufert. (2014). *Analytics and Freemium Products.* Elsevier. https://doi.org/10.1016/b978-0-12-416690-5.00002-6

Shiva Kumar Bhuram. (2025). Edge-Cloud AI for Dynamic Pricing in Automotive Aftermarkets: A Federated Reinforcement Learning Approach for Multi-Tier Ecosystems. *World Journal of Advanced Engineering Technology and Sciences.* https://doi.org/10.30574/wjaets.2025.15.3.0909.

Siddannavar, Khan, & Takalkar. (2025). Analysis of Psychological Factors Affecting Customer Lifetime Value on SaaS Platforms. *International Journal of Financial Management and Research.* https://doi.org/10.36948/ijfmr.2025.v07i04.52064.

Sonderegger. (2011). MARKET SEGMENTATION WITH NONLINEAR PRICING. **. https://doi.org/10.1111/j.1467-6451.2011.00445.x.

Subham. (2025). AI-Powered Forecasting Models for Sales and Revenue Operations. *International Journal of IoT*. https://doi.org/10.55640/ijiot-05-01-04.

Trad, & Chehab. (2024). Large Multimodal Agents for Accurate Phishing Detection with Enhanced Token Optimization and Cost Reduction. https://doi.org/10.1109/fllm63129.2024.10852444

Tucker. (2019). *Digital Data, Platforms and the Usual [Antitrust] Suspects: Network Effects, Switching Costs, Essential Facility.* https://doi.org/10.2139/ssrn.3326385

Vashishtha, Garg, & Vimal. (2022). A DATA-DRIVEN METHOD TO DYNAMIC PRICING: UNRAVELLING INVENTORY AND COMPETITOR CONTESTS WITH AI IN E-COMMERCE. *ShodhKosh: Journal of Visual and Performing Arts.* https://doi.org/10.29121/shodhkosh.v3.i2.2022.3392.

Vogt, & Bizer. (2013). *Lock-In Effects in Competitive Bidding Schemes for Payments for Ecosystem Services.* https://doi.org/10.2139/ssrn.2282039

Williamson. (2010). *Transaction Cost Economics: An Overview.* Edward Elgar Publishing. https://doi.org/10.4337/9781849806909.00007

Ye, & Zhang. (2017). Monopolistic nonlinear pricing with consumer entry. *Theoretical Economics.* https://doi.org/10.3982/te1944.

**Yin, & Qiu. (2021). *AI Technology and Online Purchase Intention:Multi-Group Analysis Based On Perceived Value.* https://doi.org/10.20944/preprints202103.0465.v1**