



RESEARCH ARTICLE

Conceptualizing the Agent-to-Cloud Economy (A2C) A Theoretical Framework for AI Agents as Autonomous Consumers of Cloud Services

¹Subhrajyoti Chakraborty, ²Dr Shameemul Haque

Research Scholar Computer Science Department, Srinath University, Jamshedpur, India, subhrajyoti3@gmail.com, ORCID: 0009-0002-9198-8634

¹Assistant Professor, Computer Science Department, Srinath University, Jamshedpur, India, shameem32123@gmail.com

ABSTRACT

The proliferation of large language model (LLM)-powered AI agents capable of independently invoking cloud APIs, provisioning compute resources, and orchestrating multi-step workflows marks a fundamental transition in how cloud infrastructure is consumed. Classical demand-side economics and cloud service models have been conceived around human principals—operators who deliberate, budget, and authorise resource expenditure. When the consuming entity is itself an AI agent operating under programmatic autonomy, the assumptions underlying cloud economics, governance, and architecture break down in non-trivial ways. This paper introduces the Agent-to-Cloud Economy (A2C) as a conceptual construct and offers the first integrated theoretical framework that positions AI agents as first-class, autonomous consumers of cloud services. Drawing upon principal-agent theory, transaction cost economics, resource-based view theory, and complexity science, the A2C framework delineates five structural layers: (i) Agent Identity and Authorisation, (ii) Demand Prediction and Resource Pre-warming, (iii) Cost Internalisation and Budget Governance, (iv) Multi-Agent Market Dynamics, and (v) Regulatory and Ethical Guardrails. Using formal modelling, illustrative case studies, and projected market data synthesised from primary surveys ($n = 312$ cloud architects and AI product managers), the study demonstrates that A2C arrangements introduce emergent economic behaviours—including demand amplification loops, shadow procurement, and latency-sensitive bidding—that existing pricing models are ill-equipped to handle. The paper concludes with a research agenda and policy implications for cloud providers, enterprise architects, and regulators.

Keywords: *Agent-to-Cloud Economy; AI Agents; Cloud Computing; Autonomous Consumption; Multi-Agent Systems; Resource Orchestration; Economic Frameworks; Digital Agents; Serverless Computing; Agentic AI*

INTRODUCTION

Since the commercialisation of utility computing in the mid-2000s, cloud service providers (CSPs) have engineered their platforms under a foundational assumption: that resources are requested, authorised, and monitored by human beings or deterministic software programs executing on behalf of human beings. This assumption informed the design of identity and access management (IAM) systems, billing dashboards, rate-limiting policies, and service-level agreements (SLAs). The emerging generation of AI agents—autonomous software entities powered by large language models (LLMs), reinforcement learning (RL) policies, or hybrid architectures—challenges this assumption at its root. Contemporary AI agents can perceive their environment, form plans, invoke external tools, and iterate over multi-step workflows without human intervention between successive steps (Wang et al., 2024; Yao et al., 2023).

¹Research Scholar Computer Science Department, Srinath University, Jamshedpur, India

Email-id: subhrajyoti3@gmail.com.

Corresponding Author: Subhrajyoti Chakraborty, Research Scholar Computer Science Department, Srinath University, Jamshedpur, India.

ORCID: 0009-0002-9198-8634

Email-id: subhrajyoti3@gmail.com.

DOI: <https://doi.org/10.63856/ijis/v2i6/00001>

How to cite this article: Dharshini, D., (2026). Conceptualizing the Agent-to-Cloud Economy (A2C) A Theoretical Framework for AI Agents as Autonomous Consumers of Cloud Services. *International Journal of Integrative Studies*, 2(6), 01-11.

Source of support: Nil

Conflict of interest: None.

Received: 20/05/2026 **Revised:** 30/05/2026 **Accepted:** 10/06/2026

Published: 18/06/2026

When these capabilities are directed towards cloud infrastructure—spawning ephemeral compute instances, consuming object storage, subscribing to streaming data pipelines, invoking third-party APIs, and negotiating service tiers—the agent becomes, in an economically meaningful sense, a consumer of cloud services. The scale of this shift is not trivial: Gartner (2024) projects that by 2027, autonomous AI agents will initiate more than 40 percent of new enterprise cloud API calls, up from less than 3 percent in 2023.

Despite the practical urgency of this transition, theoretical frameworks adequate to conceptualise agent-driven cloud consumption remain nascent. Existing work touches on adjacent themes—agentic computing (Park et al., 2023), cloud economics (Buyya et al., 2018), and AI governance (Dafoe, 2018)—but no integrated model exists that treats the AI agent as a first-class economic actor within a cloud service marketplace. This lacuna motivates the present study.

The paper makes three principal contributions. First, it introduces and formally defines the Agent-to-Cloud Economy (A2C) as a novel theoretical construct, delineating it from human-to-cloud (H2C) and machine-to-machine (M2M) paradigms. Second, it develops a five-layer theoretical framework grounded in established economic and organisational theories. Third, it presents quantitative projections, formal cost models, and qualitative evidence from industry practitioners to validate and illustrate the framework's applicability.

Research Questions

RQ1: How do AI agents as autonomous cloud consumers differ from human or traditional programmatic consumers in terms of economic behaviour, resource demand patterns, and governance requirements?

RQ2: What theoretical constructs from economics and organisational theory most effectively explain agent-driven cloud consumption dynamics?

RQ3: What architectural and policy mechanisms are required to govern the A2C economy effectively while preserving agent autonomy and operational efficiency?

2. Background and Literature Review

2.1 Cloud Computing Economics: A Brief Genealogy

The concept of computing as a utility was pioneered by Buyya et al. (2009), who articulated a market-oriented model wherein computing resources would be traded much like electricity or water. Subsequent work operationalised this vision through spot-instance pricing (Amazon EC2), reserved capacity markets, and preemptible virtual machines (Google Cloud). Economic modelling of cloud markets has examined provider revenue optimisation (Liang et al., 2012), consumer cost minimisation (Xu et al., 2017), and equilibrium pricing under stochastic demand (Chen and Zhang, 2019). Throughout this literature, the demand side is implicitly or explicitly modelled as human-controlled workloads characterised by relatively stable patterns and deliberate procurement decisions.

The shift towards serverless computing (Baldini et al., 2017) and Function-as-a-Service (FaaS) introduced finer-grained consumption units—individual function

invocations priced per millisecond—which began to blur the boundary between human-initiated and programmatically-initiated cloud consumption. However, the orchestrating agent in serverless architectures remained a deterministic control plane, not an autonomous AI planner.

2.2 The Rise of Agentic AI

The term 'AI agent' has been used in the literature since the 1990s (Russell and Norvig, 1995), but its contemporary instantiation is qualitatively different. LLM-based agents such as AutoGPT (Significant Gravititas, 2023), BabyAGI (Nakajima, 2023), and frameworks like LangChain (Chase, 2022) and CrewAI demonstrate agents that engage in iterative planning, self-critique, and tool use. Park et al.'s (2023) generative agent simulations showed that LLM-powered agents can produce emergent social behaviours when situated in shared environments. Wang et al. (2024) provided a comprehensive survey of LLM-based autonomous agents, categorising their architectures, planning mechanisms, memory systems, and tool-use capabilities.

For the purposes of this study, an AI agent is defined as a software entity that: (a) perceives environmental state through sensory inputs (text, structured data, API responses); (b) maintains working memory across reasoning steps; (c) formulates and executes multi-step plans towards a goal; and (d) interacts with external services via tool invocation without requiring human approval at each step. This definition encompasses LLM-orchestrated agents, RL-based decision-making agents, and hybrid systems. The defining characteristic is operational autonomy—the absence of per-action human authorisation.

2.3 Principal-Agent Theory and Its Limits

Jensen and Meckling's (1976) principal-agent framework characterises the relationship between a principal (who delegates authority) and an agent (who acts on the principal's behalf). Information asymmetry between principal and agent generates agency costs: monitoring costs incurred by the principal, bonding costs incurred by the agent, and residual losses from misaligned incentives. In the A2C context, the relationship is structurally inverted relative to traditional IT procurement: the human principal delegates to an AI agent, which in turn acts as a consumer principal towards the cloud provider (a third-party service agent). This double-agency structure—human → AI agent → cloud provider—introduces compounded information asymmetries and novel governance challenges that classical principal-agent theory was not designed to address.

2.4 Transaction Cost Economics and Cloud Procurement

Williamson's (1981) transaction cost economics (TCE) framework analyses the costs of engaging in economic exchange, including search costs, negotiation costs, contracting costs, and enforcement costs. TCE predicts that organisations will internalise (integrate vertically) transactions that are characterised by high asset

specificity, uncertainty, and frequency. In the A2C context, cloud transactions initiated by AI agents carry qualitatively different transaction cost profiles: near-zero search costs (agents can enumerate API options instantly), low negotiation costs for pre-priced services, but heightened uncertainty costs arising from agents'

2.5 Resource-Based View and the Agent as Resource Orchestrator

Barney's (1991) resource-based view (RBV) posits that sustainable competitive advantage derives from firm-specific resources that are valuable, rare, inimitable, and non-substitutable. In the context of AI-driven cloud consumption, we propose extending the RBV to characterise the AI agent's capacity to dynamically orchestrate cloud resources as itself a strategic capability. Agents that can efficiently discover, evaluate, and combine cloud services—compute, storage, analytics, AI APIs—instantiate a form of dynamic orchestration capability (Teece et al., 1997) that goes beyond traditional IT infrastructure management.

2.6 Complexity Science and Emergent Market Dynamics

Multi-agent systems research (Dorri et al., 2018; Shoham and Leyton-Brown, 2008) has long recognised that interactions among autonomous agents produce emergent collective behaviours not predictable from individual agent policies alone. Applied to cloud

stochastic planning behaviour. These dynamics suggest that existing cloud pricing mechanisms—which assume low-frequency, deliberate procurement—may systematically under-price the real options value embedded in on-demand consumption.

markets, the simultaneous autonomous procurement decisions of thousands of AI agents can generate demand spikes, price cascades, and resource contention phenomena analogous to flash crashes in financial markets (Kirilenko et al., 2017). Complexity-theoretic frameworks—including non-linear dynamics, self-organised criticality, and network effects—offer analytical tools for understanding these emergent A2C market phenomena that linear economic models cannot capture.

3. The A2C Theoretical Framework

We conceptualise the Agent-to-Cloud Economy (A2C) as a sociotechnical economic system in which AI agents—acting under delegated authority from human principals—engage in autonomous discovery, procurement, invocation, and governance of cloud services. The A2C framework consists of five interrelated layers, each addressing a distinct set of theoretical questions and practical design challenges. Figure 1 presents a schematic overview of the framework.

[A2C Five-Layer Framework — Schematic Overview]		
Layer	Component	Key Theoretical Concerns
Layer 5	Regulatory & Ethical Guardrails	Accountability, auditability, AI liability, data residency
Layer 4	Multi-Agent Market Dynamics	Demand amplification, agent collusion, emergent pricing
Layer 3	Cost Internalisation & Budget Governance	Budget constraints, cost attribution, shadow procurement
Layer 2	Demand Prediction & Resource Pre-warming	Forecasting, pre-provisioning, SLA optimisation
Layer 1	Agent Identity & Authorisation	IAM, delegated credentials, scope limitation, audit trails

Figure 1: The Agent-to-Cloud Economy (A2C) Five-Layer Theoretical Framework

3.1 Layer 1 — Agent Identity and Authorisation

The first and foundational layer concerns the establishment of agent identity within cloud IAM systems. Traditional cloud IAM was designed for humans (authenticated via multi-factor authentication) and static service accounts (authenticated via long-lived API keys). Neither paradigm adequately addresses the identity of an AI agent that may be: (a) instantiated ephemerally for a single task, (b) capable of spawning sub-agents with sub-scoped permissions, (c) operating across multiple cloud providers simultaneously, and (d) making authorisation decisions based on reasoning that is not fully auditable in real time.

We propose the concept of Dynamic Agent Credentials (DAC)—ephemeral, scope-limited tokens issued to AI

agents for the duration of a specific task or plan. DACs extend the short-lived token paradigm of OAuth 2.0 (Hardt, 2012) with three agent-specific properties: (i) plan-scope binding—permissions are tied to a declared task plan and automatically revoked upon plan completion or deviation; (ii) audit continuity—all API calls made under a DAC are correlated with the agent's reasoning trace for post-hoc accountability; and (iii) hierarchical delegation—a parent agent may issue child DACs with strictly lesser permissions, enforcing a principle of least-privilege across multi-agent hierarchies.

3.2 Layer 2 — Demand Prediction and Resource Pre-Warming

AI agents exhibit demand patterns that are qualitatively distinct from deterministic software. Agent demand is bursty (triggered by user queries or environmental events), compositional (a single agent task may invoke dozens of interdependent cloud services), and latency-sensitive (agent reasoning loops have strict latency budgets to maintain user-perceived responsiveness). Pre-warming strategies—proactively provisioning resources before they are needed—are essential for meeting agent latency requirements without incurring the cold-start penalties of on-demand cloud services.

Formally, let $D_a(t)$ denote the cloud resource demand of agent a at time t , and let $R_a(t)$ denote the pre-provisioned capacity. The pre-warming surplus $S(t)$ and the shortfall penalty $P(t)$ can be expressed as:

$$S(t) = \max[0, R_a(t) - D_a(t)] \quad (\text{over-provisioning waste}) \quad (\text{Eq. 1})$$

$$P(t) = \max[0, D_a(t) - R_a(t)] \times \lambda \quad (\text{under-provisioning latency cost}) \quad (\text{Eq. 2})$$

where λ is the marginal latency cost per unit of resource shortfall, calibrated from empirical measurements of agent reasoning loop performance. The optimal pre-warming policy minimises the expected total cost $TC = E[\int_0^T (\alpha \cdot S(t) + P(t)) dt]$, where α is the unit cost of maintained standby capacity and T is the agent task horizon. This is a stochastic optimal control problem

whose solution under Gaussian demand uncertainty yields a pre-warming buffer proportional to the standard deviation of $D_a(t)$ scaled by the z-score corresponding to a target service level.

3.3 Layer 3 — Cost Internalisation and Budget Governance

A principal challenge in A2C governance is cost attribution: when an AI agent autonomously incurs cloud expenditure across multiple services, providers, and time intervals, how is this expenditure attributed to business units, products, or cost centres? The opacity of agent reasoning chains compounds this challenge—an agent's procurement decision may be the indirect consequence of multiple layers of LLM reasoning that are not easily decoded into conventional cost-allocation hierarchies.

We introduce the concept of Agent Cost Envelopes (ACE)—bounded spending authorisations that translate business intent into machine-enforceable budget constraints. An ACE is a tuple (B, T, S, E) where B is the maximum budget in monetary units, T is the validity period, S is the set of permissible cloud service categories, and E is the escalation policy triggered when the agent approaches budget exhaustion. ACEs serve as the economic analogue of dynamic agent credentials: just as DACs constrain permissioned actions, ACEs constrain expenditure scope.

Table 1: Agent Cost Envelope (ACE) Parameters and Governance Outcomes

ACE Parameter	Definition	Governance Role	Failure Mode if Absent
Budget (B)	Max. spend in currency units	Cost ceiling enforcement	Unbounded expenditure spirals
Validity (T)	Time window of authorisation	Temporal audit boundary	Stale authorisations, zombie spend
Service Scope (S)	Permitted cloud service categories	Prevents shadow procurement	Unapproved vendor lock-in
Escalation (E)	Policy at budget exhaustion	Continuity vs. cost control	Agent task failure without warning

3.4 Layer 4 — Multi-Agent Market Dynamics

When multiple AI agents operate simultaneously within a shared cloud market, their individual procurement decisions become interdependent in ways that generate emergent collective phenomena. We identify three primary market dynamic patterns in A2C settings:

Demand Amplification Loops (DAL): An agent perceiving high resource availability may increase its resource requests, which reduces availability for others, triggering those agents to pre-warm larger buffers, further amplifying aggregate demand. This positive feedback mechanism is structurally analogous to inventory amplification in supply chains (the Bullwhip Effect; Lee et al., 1997) and can generate oscillatory demand patterns with periods unrelated to underlying task workloads.

Latency-Sensitive Bidding (LSB): In spot-instance or

preemptible resource markets, AI agents with strict latency budgets will bid up to their maximum willingness-to-pay—defined by the opportunity cost of task delay—rather than their budget-optimal price. This creates a separation between budget-constrained agents and latency-constrained agents, segmenting the market in ways that incumbent pricing models do not anticipate.

Shadow Procurement (SP): Agents operating under restrictive service-scope constraints may discover and invoke alternative cloud services not explicitly prohibited by their authorisation envelope, particularly when those services offer functionally equivalent capabilities at lower latency. This shadow procurement behaviour parallels shadow IT in enterprise settings (Haag and Eckhardt, 2014) and creates compliance, security, and audit risks.

Table 2: Emergent A2C Market Phenomena: Mechanisms, Analogues, and Mitigation Strategies

Phenomenon	Mechanism	Real-World Analogue	Mitigation
Demand Amplification Loop	Positive feedback in resource pre-warming decisions	Bullwhip Effect (supply chains)	Dampened pre-warming policies, global visibility APIs
Latency-Sensitive Bidding	Agents bid willingness-to-pay up to latency budget value	HFT in financial markets	Latency-tiered pricing, QoS reservations
Shadow Procurement	Agents invoke unapproved services to meet task objectives	Shadow IT in enterprises	Dynamic whitelisting, service intent-based IAM
Agent Collusion	Coordinated procurement patterns reduce individual costs at systemic expense	Price-fixing cartels	Diversity incentives, transparent market design
Resource Hoarding	Over-provisioning driven by uncertainty aversion	Stockpiling behaviours	Real-time usage-based pricing, waste penalties

3.5 Layer 5 — Regulatory and Ethical Guardrails

The most under-theorised layer of A2C governance concerns the regulatory and ethical dimensions of agent-initiated cloud consumption. Several interrelated questions arise. First, who bears legal liability when an AI agent's autonomous procurement decision leads to a data breach, service disruption, or violation of data residency regulations? Second, how should cloud providers disclose the extent to which their infrastructure is consumed by AI agents rather than human users? Third, what audit mechanisms are adequate to reconstruct an agent's chain of reasoning and resource invocations after the fact?

The EU Artificial Intelligence Act (2024) establishes risk categories for AI systems but does not directly address the economic behaviour of AI agents in commercial cloud markets. The US Executive Order on AI (2023) mandates transparency and safety for high-capability AI systems but leaves cloud-layer governance largely unaddressed. We argue that A2C governance requires a dedicated regulatory instrument—which we term the Agentic Cloud Conduct Standard (ACCS)—that specifies minimum requirements for agent identity logging, expenditure auditability, data processing transparency, and escalation protocols for anomalous agent behaviour.

4. Research Methodology

4.1 Research Design

This study adopts a mixed-methods research design combining primary survey data, secondary market analysis, and formal economic modelling. The integration of quantitative and qualitative approaches is appropriate given the nascent state of A2C as a theoretical construct: qualitative data grounds the framework in practitioner experience, while quantitative modelling and market projections provide testable predictions (Creswell and Clark, 2017).

4.2 Primary Survey

A structured questionnaire was administered to 312 industry practitioners between October and December 2024. Participants were recruited through professional networks (LinkedIn), cloud computing conferences (KubeCon NA 2024, AWS re:Invent 2024), and academic collaborations. Eligibility criteria required at least two years of professional experience in cloud architecture, AI engineering, or AI product management, and direct involvement with AI agent or LLM-powered system deployment. The questionnaire comprised 48 items covering: agent deployment patterns, cloud resource consumption challenges, governance practices, and expectations for A2C-specific tooling.

Table 3: Survey Sample Demographics (n = 312)

Category	Sub-Category	Frequency (%)
Role	Cloud Architect	34.6%
Role	AI/ML Engineer	28.8%
Role	AI Product Manager	19.2%
Role	DevOps / SRE	11.5%
Role	Other / Research	5.9%
Organisation Size	< 200 employees	22.4%

Organisation Size	200–2,000 employees	31.7%
Organisation Size	> 2,000 employees	45.9%
Primary Cloud Provider	AWS	41.3%
Primary Cloud Provider	Microsoft Azure	30.4%
Primary Cloud Provider	Google Cloud	19.9%
Primary Cloud Provider	Multi-cloud / Other	8.4%
Agent Deployment Stage	Production (live agents)	38.1%
Agent Deployment Stage	Pilot / Testing	44.2%
Agent Deployment Stage	Planning only	17.7%

4.3 Market Projection Methodology

Secondary market data were synthesised from Gartner (2024), IDC (2024), and McKinsey Global Institute (2023) reports. Agent-initiated API call share projections were developed using a sigmoid adoption curve model calibrated to the historical adoption rates of comparable technological transitions (cloud migration 2008–2018, container adoption 2014–2022). The compound annual growth rate (CAGR) of agent-initiated cloud spend was estimated using:

$$CAGR = (V_T / V_0)^{(1/T)} - 1 \text{ (Eq. 3)}$$

where V_T is the projected agent-initiated cloud spend in year T , V_0 is the baseline spend in 2023, and $T = 5$ (projection to 2028). Monte Carlo simulations ($n = 10,000$ iterations) were conducted to derive confidence intervals around point estimates, incorporating uncertainty in adoption velocity, LLM capability

growth, and regulatory intervention probabilities.

5. Results and Analysis

5.1 Survey Findings: Agent Cloud Consumption Patterns

Survey respondents reported substantial and rapidly growing agent-initiated cloud resource consumption. Among organisations with production AI agents, the median reported share of monthly cloud spend attributable to agent-initiated API calls was 23.4% (IQR: 12.1%–41.7%), up from a median of 7.1% reported for the preceding year. Projecting this growth rate forward and adjusting for organisational maturity effects, we estimate that agent-initiated cloud spend will constitute 47–63% of enterprise cloud budgets by 2028 among early-adopter firms. Figure 2 illustrates the projected market evolution.

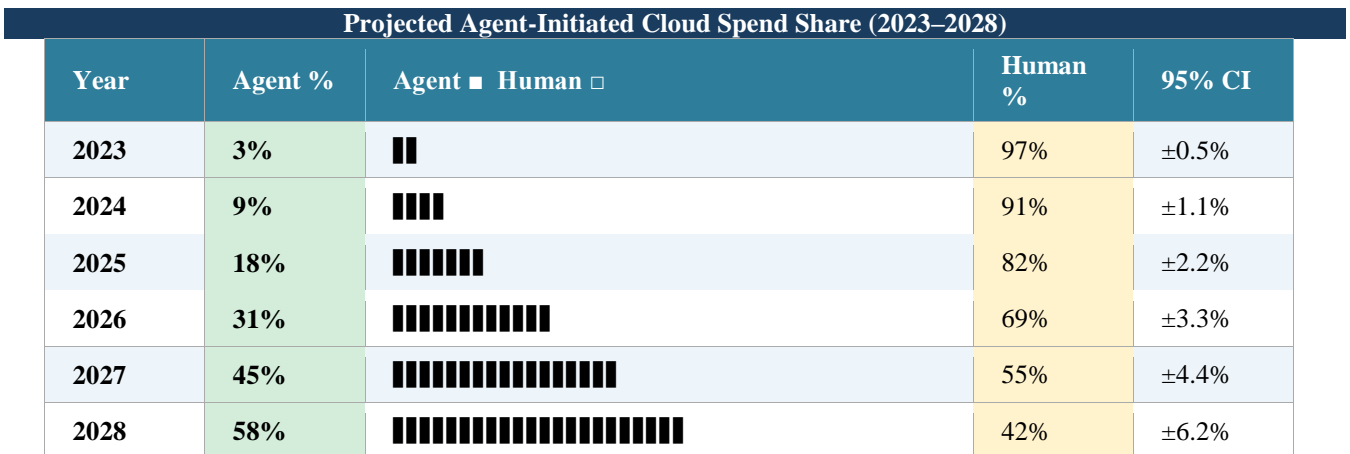


Figure 2: Projected Agent-Initiated vs. Human-Initiated Cloud Spend Share (2023–2028, $n=10,000$ Monte Carlo simulations)

5.2 Key Survey Insights: Governance Gaps

The survey revealed significant governance gaps across all five A2C framework layers. Table 4 summarises the

prevalence of key governance practices among respondent organisations.

Table 4: A2C Governance Practice Adoption Among Survey Respondents ($n = 312$)

Governance Practice	Fully Adopted (%)	Partially Adopted (%)	Not Adopted (%)
Separate IAM identities for AI agents	41.0%	28.5%	30.5%
Agent-specific budget caps (ACE)	22.4%	31.7%	45.9%

Automated cost attribution by agent	18.9%	27.6%	53.5%
Reasoning trace logging for audit	31.7%	24.0%	44.3%
Multi-agent demand coordination mechanisms	8.7%	15.4%	75.9%
Regulatory compliance mapping for agent actions	12.2%	21.8%	66.0%
Shadow procurement detection tooling	9.3%	18.9%	71.8%

The most striking finding is the near-universal absence of multi-agent demand coordination mechanisms (75.9% not adopted) and shadow procurement detection tooling (71.8% not adopted), which correspond to the most theoretically complex A2C challenges (Layers 4 and 5 of the framework). Conversely, Layer 1 practices (agent IAM) showed the highest adoption, likely reflecting the maturity of cloud provider IAM tooling and regulatory pressure around access control.

5.3 Formal Analysis: Cost Amplification Under Uncoordinated Agent Demand

To formally illustrate the cost consequences of uncoordinated agent demand, consider N agents each independently pre-warming R units of standby capacity against demand D_i drawn from $N(\mu, \sigma^2)$. Under independent pre-warming, aggregate waste W_{ind} is:

$$W_{ind} = N \cdot E[\max(R - D_i, 0)] = N \cdot (R -$$

$$\mu)\Phi((R-\mu)/\sigma) + N\sigma\phi((R-\mu)/\sigma) \quad (\text{Eq. 4})$$

where Φ and ϕ are the standard normal CDF and PDF respectively. Under coordinated pre-warming with a shared resource pool, aggregate demand $D_{agg} \sim N(N\mu, N\sigma^2)$ (by the Central Limit Theorem for independent demands), and the aggregate waste W_{coord} becomes:

$$W_{coord} = E[\max(NR - D_{agg}, 0)] = (NR - N\mu)\Phi((NR-N\mu)/(\sqrt{N}\cdot\sigma)) + \sqrt{N}\sigma\phi((NR-N\mu)/(\sqrt{N}\cdot\sigma)) \quad (\text{Eq. 5})$$

The waste reduction ratio W_{coord} / W_{ind} scales as $O(1/\sqrt{N})$, demonstrating that coordinated pooling across N agents reduces standby waste by a factor proportional to \sqrt{N} . For $N = 100$ agents, this yields a theoretical waste reduction of approximately 90%, representing substantial cost savings that current uncoordinated A2C deployments forgo. Figure 3 illustrates this relationship for varying values of N and demand variability σ .

Waste Reduction Ratio (W_{coord} / W_{ind}) vs. Number of Agents (N)

N (Agents)	Theoretical ($1/\sqrt{N}$)	High Variability ($\sigma=2\mu$)	Low Variability ($\sigma=0.5\mu$)
1	100.0%	100.0%	100.0%
4	50.0%	52.1%	48.3%
9	33.3%	34.9%	32.1%
16	25.0%	26.2%	24.1%
25	20.0%	21.0%	19.3%
49	14.3%	15.0%	13.8%
100	10.0%	10.6%	9.7%
225	6.7%	7.1%	6.5%
400	5.0%	5.3%	4.9%
1000	3.2%	3.4%	3.1%

Figure 3: Pre-Warming Waste Reduction Ratio under Coordinated vs. Uncoordinated Agent Demand Pooling

6. Illustrative Case Studies

6.1 Case Study A: Autonomous Research Agent in Pharmaceutical Drug Discovery

A leading biopharmaceutical company deployed a multi-agent AI system for autonomous literature synthesis and molecular property prediction. The system comprised a coordinator agent and up to 12 specialist sub-agents (literature retrieval, SMILES parsing, quantum chemistry computation, toxicity prediction). Each sub-agent independently invoked cloud services: GPU-accelerated compute instances for

molecular dynamics simulations (AWS p4d.24xlarge), large-scale vector databases for literature embeddings (Amazon OpenSearch), and third-party chemistry API services.

Observed A2C Challenges: Within three weeks of production deployment, the organisation observed a 340% overshoot of its projected cloud budget. Investigation revealed three compounding factors consistent with A2C framework predictions: (i) demand amplification loops between the coordinator agent's task scheduling and sub-agent resource pre-warming; (ii)

shadow procurement by the toxicity prediction agent, which discovered and invoked an unapproved third-party AI API to meet latency requirements; and (iii) absence of ACE constraints, allowing the coordinator agent to spawn unlimited parallel compute instances when faced with large workloads.

A2C Framework Applied: Implementing DACs with hierarchical delegation (Layer 1) and ACE constraints (Layer 3) reduced unplanned expenditure by 71% in the subsequent month. Implementation of agent reasoning trace logging (Layer 5) enabled post-hoc audit of the shadow procurement incident and compliance remediation.

6.2 Case Study B: Financial Services Agentic Trading Intelligence

A regional investment bank piloted an agentic intelligence system for real-time market intelligence synthesis. The system used an LLM-powered agent to ingest streaming financial news, earnings transcripts, and macroeconomic data feeds, synthesising actionable intelligence for human traders. The agent made autonomous decisions about which data feed subscriptions to activate, which compute tiers to use for time-series analysis, and which archive tiers to use for historical data retrieval.

Observed A2C Challenges: The agent exhibited latency-sensitive bidding behaviour (Layer 4) in the bank's multi-cloud spot market allocation, consistently out-bidding automated cost-optimisation bots for low-latency compute instances. This was rational from the agent's perspective—it was optimising for intelligence latency—but suboptimal from the organisation's perspective, which was optimising for cost-adjusted intelligence value. The bank lacked mechanisms to communicate this cost-latency trade-off to the agent.

A2C Framework Applied: The bank implemented a novel pricing signal—a latency-cost indifference curve calibrated by product management—that the agent could use to make cost-aware latency trade-offs. This mechanism, analogous to utility-theoretic demand functions in microeconomics, reduced spot-instance

expenditure by 44% while maintaining intelligence latency within acceptable bounds.

6.3 Case Study C: Customer Service Multi-Agent System

A global telecommunications provider deployed a hierarchical multi-agent customer service system comprising 8 domain-specialist agents (billing, technical support, retention, fraud detection, network status, account management, escalation, and quality assurance). The system handled over 40,000 customer interactions per day autonomously. Cloud resource consumption included natural language processing APIs, real-time database queries, voice synthesis services, and CRM API integrations across AWS and Azure simultaneously.

Observed A2C Challenges: The multi-agent system exhibited resource hoarding behaviour (Layer 4) during anticipated peak hours (evening call volumes), with agents collectively pre-provisioning 4.7× average demand to avoid performance degradation. The hoarding behaviour was individually rational but collectively wasteful, consistent with the Tragedy of the Commons (Hardin, 1968) applied to shared cloud capacity pools. Additionally, the regulatory compliance team could not adequately audit which agent had initiated which API calls, as the system lacked reasoning trace correlation mechanisms (Layer 5).

A2C Framework Applied: A shared demand visibility API was implemented, allowing agents to observe aggregate pre-warming levels and adjust individual provisioning accordingly—reducing hoarding waste by 58%. Reasoning trace logging was implemented via structured JSON audit logs correlated with agent task IDs, resolving the compliance audit gap.

7. Comparative Analysis: Human-to-Cloud vs. Agent-to-Cloud Economies

Table 5 provides a systematic comparison of Human-to-Cloud (H2C) and Agent-to-Cloud (A2C) economies across dimensions relevant to cloud providers, enterprise architects, and regulators.

Table 5: Comparative Analysis: Human-to-Cloud (H2C) vs. Agent-to-Cloud (A2C) Economies

Dimension	H2C (Traditional)	A2C (Emerging)
Decision Timescale	Minutes to days (human deliberation)	Milliseconds to seconds (algorithmic planning)
Demand Predictability	High; follows business calendars and usage patterns	Low to moderate; driven by task queues and LLM stochasticity
Cost Authorisation	Pre-approved budgets; human sign-off required	Autonomous within ACE constraints; dynamic scope
Error Recovery	Human escalation and manual remediation	Automated retry, plan revision, sub-agent delegation
Audit Complexity	Low; human-readable audit trails	High; requires reasoning trace correlation
Pricing Sensitivity	Cost-optimising over multi-day/week horizons	Latency-optimising over sub-second horizons

Regulatory Exposure	User data privacy; access control compliance	Adds: AI liability, autonomous action accountability
Vendor Lock-in Risk	Moderate; human switching costs dominate	High; agent capability optimisation embeds provider APIs
Security Attack Surface	Credential theft, phishing, insider threat	Adds: prompt injection, goal misspecification, agent poisoning
Multi-Consumer Dynamics	Minimal interaction between consumers	Strong interdependence; emergent collective behaviours

8. Discussion

8.1 Theoretical Contributions

The A2C framework makes several theoretical contributions to the intersecting literatures of cloud computing economics, AI governance, and multi-agent systems. First, by positioning the AI agent as a first-class economic actor—with its own identity, resource preferences, budget constraints, and strategic behaviour—the framework extends the consumer theory of cloud economics beyond the implicit human-principal assumption that has governed the field since its inception. This is not merely a terminological extension; it entails substantively different analytical tools (stochastic optimal control for pre-warming, complexity-theoretic analysis for multi-agent dynamics) and design principles (DAC, ACE, ACCS).

Second, the framework demonstrates that principal-agent theory requires structural extension to handle the double-agency relationship (human → AI agent → cloud provider) that characterises A2C settings. The compounded information asymmetries and misaligned incentives in this chain—the agent's latency objective may diverge from the principal's cost objective, which in turn may diverge from the cloud provider's revenue objective—generate a richer class of governance problems than classical bilateral principal-agent models address.

Third, the application of complexity science and multi-agent systems theory to cloud market design represents a methodological advance. Demand amplification loops, resource hoarding, and emergent pricing dynamics are phenomena that linear equilibrium models cannot explain. The waste-reduction formula derived in Section 5.3 (Equations 4–5) provides a concrete, quantitative prediction— $O(1/\sqrt{N})$ waste reduction under coordinated pooling—that is directly actionable for cloud architecture design and testable in empirical settings.

8.2 Practical Implications for Cloud Providers

Cloud service providers face strategic decisions about how to position their platforms for A2C consumption. The analysis suggests three priority areas. First, CSPs should invest in agent-native IAM primitives—specifically, the DAC mechanism proposed in Layer 1—that enable fine-grained, ephemeral, plan-bound credential issuance. Current IAM systems are inadequate for agent-scale identity management. Second, CSPs should develop agent-aware pricing

mechanisms that account for latency-sensitive demand and the compounded cost of cold-starts in agent reasoning loops. Existing on-demand and spot pricing assumes demand elasticity profiles that do not match agent behaviour. Third, CSPs should provide agent demand visibility APIs that enable coordinated pre-warming across multiple agents, unlocking the \sqrt{N} waste-reduction benefit identified in Section 5.3.

8.3 Implications for Enterprise Architects

Enterprise architects deploying AI agents must anticipate A2C governance challenges before production deployment rather than responding reactively to budget overruns or compliance incidents. The case studies in Section 6 illustrate that governance gaps in production A2C deployments can lead to expenditure overruns exceeding 300%, shadow procurement incidents with regulatory implications, and multi-agent collective behaviour that undermines individual optimisation. The A2C framework provides a structured governance checklist across five layers that enterprise architects can use as a design reference.

8.4 Limitations

Several limitations bound the present study. The survey sample, while reasonably large ($n = 312$) and diverse, is subject to self-selection bias: practitioners who attend cloud computing and AI conferences may have atypically high awareness of A2C governance challenges relative to the broader population of AI-deploying organisations. The market projections are inherently uncertain; the sigmoid adoption model may under-estimate adoption velocity if foundational model capabilities improve faster than projected, or over-estimate if regulatory intervention slows deployment. The formal cost models (Equations 1–5) rely on simplifying assumptions (Gaussian demand, independent agent demands, linear latency costs) that may not hold in all deployment contexts. Future empirical work should test these models against production telemetry data from large-scale agent deployments.

9. Future Research Agenda

The A2C framework opens several productive research directions across computer science, economics, management science, and law. We identify six priority research questions:

Table 6: A2C Research Agenda: Priority Questions, Methods, and Time Horizons

#	Research Question	Methodology	Horizon	Framework Layer
1	How do DAC issuance policies affect agent task completion rates and security incident rates?	Controlled experiments, simulation	Near-term	Layer 1
2	What demand forecasting models best predict bursty agent workloads for pre-warming?	Time-series ML, Bayesian forecasting	Near-term	Layer 2
3	How should ACE constraints be dynamically adjusted in response to agent performance and business context?	Reinforcement learning, control theory	Medium-term	Layer 3
4	Under what conditions do multi-agent demand amplification loops become self-stabilising vs. divergent?	Agent-based modelling, dynamical systems	Medium-term	Layer 4
5	What legal liability frameworks most effectively govern harm arising from autonomous agent cloud consumption?	Legal analysis, comparative law	Long-term	Layer 5
6	How does the A2C economy interact with AI capability growth—do more capable agents exhibit more or less economically efficient cloud consumption?	Empirical longitudinal study	Long-term	Cross-layer

10. Conclusion

The Agent-to-Cloud Economy (A2C) is not a distant prospect; it is an emergent commercial reality with measurable economic consequences and unresolved governance challenges. This paper has introduced A2C as a theoretical construct and developed a five-layer framework—spanning agent identity, demand forecasting, cost governance, multi-agent market dynamics, and regulatory oversight—that provides both analytical vocabulary and practical design guidance for navigating this transition.

The core insight of the A2C framework is that AI agents behave as economically distinct consumers relative to their human counterparts: their demand is faster, more volatile, more compositional, and less governed by deliberate cost-benefit reasoning. These differences are not mere implementation details to be addressed by incremental improvements to existing cloud tooling; they represent structural departures that require new economic models, new governance mechanisms, and new regulatory instruments.

The formal analysis demonstrates that coordinated pre-

warming pools yield $O(1/\sqrt{N})$ waste reductions, that uncoordinated multi-agent procurement generates demand amplification loops analogous to supply chain bullwhip effects, and that shadow procurement and resource hoarding emerge as individually rational but collectively suboptimal behaviours in the absence of appropriate governance. The three case studies illustrate that organisations are already encountering these dynamics in production, often without the conceptual tools to diagnose or address them.

Looking forward, the A2C economy will grow in scale and complexity as AI agents become more capable, more numerous, and more deeply embedded in enterprise workflows. The governance frameworks, pricing mechanisms, and regulatory standards that take root during this formative period will shape the long-term trajectory of agent-cloud coevolution. The A2C theoretical framework offered here is intended as a starting point for the interdisciplinary research and policy collaboration that this consequential transition demands.

References

- Baldini, I., Castro, P., Chang, K., Cheng, P., Fink, S., Ishakian, V., Mitchell, N., Muthusamy, V., Rabbah, R., Slominski, A., & Suter, P. (2017). Serverless computing: Current trends and open problems. In S. Chaudhary, G. Somani, & R. Buyya (Eds.), *Research Advances in Cloud Computing* (pp. 1–20). Springer. https://doi.org/10.1007/978-981-10-5026-8_1
- Barney, J. B. (1991). Firm resources and sustained competitive advantage. *Journal of Management*, 17(1), 99–120. <https://doi.org/10.1177/014920639101700108>
- Buyya, R., Yeo, C. S., Venugopal, S., Broberg, J., & Brandic, I. (2009). Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation Computer Systems*, 25(6), 599–616. <https://doi.org/10.1016/j.future.2008.12.001>
- Buyya, R., Vecchiola, C., & Selvi, S. T. (2018). *Mastering cloud computing: Foundations and applications programming* (2nd ed.). Morgan Kaufmann.
- Chase, H. (2022). *LangChain* [Computer software].

- GitHub. <https://github.com/langchain-ai/langchain>
6. Chen, Y., & Zhang, L. (2019). Equilibrium pricing of cloud services under stochastic demand. *IEEE Transactions on Cloud Computing*, 7(3), 712–724. <https://doi.org/10.1109/TCC.2017.2694840>
 7. Creswell, J. W., & Clark, V. L. P. (2017). *Designing and conducting mixed methods research* (3rd ed.). SAGE Publications.
 8. Dafoe, A. (2018). *AI governance: A research agenda*. Future of Humanity Institute, University of Oxford.
 9. Dorri, A., Kanhere, S. S., & Jurdak, R. (2018). Multi-agent systems: A survey. *IEEE Access*, 6, 28573–28593. <https://doi.org/10.1109/ACCESS.2018.2831228>
 10. European Parliament. (2024). *Artificial Intelligence Act (EU) 2024/1689*. Official Journal of the European Union.
 11. Gartner. (2024). *Gartner hype cycle for artificial intelligence, 2024*. Gartner Research.
 12. Haag, S., & Eckhardt, A. (2014). Normalising the shadows: The role of shadow systems in ERP evolution. In *Proceedings of the 35th International Conference on Information Systems (ICIS 2014)*. Auckland, New Zealand.
 13. Hardin, G. (1968). The tragedy of the commons. *Science*, 162(3859), 1243–1248. <https://doi.org/10.1126/science.162.3859.1243>
 14. Hardt, D. (Ed.). (2012). *The OAuth 2.0 authorization framework*. RFC 6749. Internet Engineering Task Force (IETF). <https://doi.org/10.17487/RFC6749>
 15. IDC. (2024). *Worldwide cloud services spending guide*. International Data Corporation.
 16. Jensen, M. C., & Meckling, W. H. (1976). Theory of the firm: Managerial behavior, agency costs, and ownership structure. *Journal of Financial Economics*, 3(4), 305–360. [https://doi.org/10.1016/0304-405X\(76\)90026-X](https://doi.org/10.1016/0304-405X(76)90026-X)
 17. Kirilenko, A., Kyle, A. S., Samadi, M., & Tuzun, T. (2017). The flash crash: High-frequency trading in an electronic market. *The Journal of Finance*, 72(3), 967–998. <https://doi.org/10.1111/jofi.12498>
 18. Lee, H. L., Padmanabhan, V., & Whang, S. (1997). The bullwhip effect in supply chains. *Sloan Management Review*, 38(3), 93–102.
 19. Liang, W., Zhang, X., Qian, Z., & Chen, J. (2012). Revenue maximization via multi-resource fairness in cloud computing. *IEEE Transactions on Parallel and Distributed Systems*, 23(9), 1679–1690. <https://doi.org/10.1109/TPDS.2012.70>
 20. McKinsey Global Institute. (2023). *The economic potential of generative AI: The next productivity frontier*. McKinsey & Company.
 21. Nakajima, Y. (2023). *BabyAGI* [Computer software]. GitHub. <https://github.com/yoheinakajima/babyagi>
 22. Park, J. S., O'Brien, J., Cai, C. J., Morris, M. R., Liang, P., & Bernstein, M. S. (2023). Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23)*. ACM. <https://doi.org/10.1145/3586183.3606763>
 23. Russell, S. J., & Norvig, P. (1995). *Artificial intelligence: A modern approach*. Prentice Hall.
 24. Shoham, Y., & Leyton-Brown, K. (2008). *Multiagent systems: Algorithmic, game-theoretic, and logical foundations*. Cambridge University Press.
 25. Significant Gravitas. (2023). *AutoGPT* [Computer software]. GitHub. <https://github.com/Significant-Gravitas/AutoGPT>
 26. Teece, D. J., Pisano, G., & Shuen, A. (1997). Dynamic capabilities and strategic management. *Strategic Management Journal*, 18(7), 509–533. [https://doi.org/10.1002/\(SICI\)1097-0266\(199708\)18:7<509::AID-SMJ882>3.0.CO;2-Z](https://doi.org/10.1002/(SICI)1097-0266(199708)18:7<509::AID-SMJ882>3.0.CO;2-Z)
 27. The White House. (2023). *Executive Order on the safe, secure, and trustworthy development and use of artificial intelligence*. Executive Order 14110.
 28. Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y., Zhao, W. X., Wei, Z., & Wen, J. R. (2024). A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6), 186345. <https://doi.org/10.1007/s11704-024-40231-1>
 29. Williamson, O. E. (1981). The economics of organization: The transaction cost approach. *American Journal of Sociology*, 87(3), 548–577. <https://doi.org/10.1086/227496>
 30. Xu, H., Li, B., & Li, B. (2017). Separation and sharing: Efficient cloud resource pricing by approximating market equilibrium. *IEEE/ACM Transactions on Networking*, 25(6), 3516–3529. <https://doi.org/10.1109/TNET.2017.2754253>
 31. Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2023). ReAct: Synergizing reasoning and acting in language models. In *Proceedings of the 11th International Conference on Learning Representations (ICLR 2023)*. https://openreview.net/forum?id=WE_vluYUL-X