**Research Article**

# Becoming a Cloud-Native Architect in the LLM Era: Skills, Tools, and Mindset for the Next Generation

**Naveen Kumar Jayakumar\***

Independent Researcher, USA
\* **Corresponding Author Email:** naveenkumar.jayakumar@gmail.com- **ORCID:** 0000-0002-0047-9950

**Abstract:**

Cloud-native approaches have reshaped the software architect's role from producing static designs to guiding decisions in complex distributed systems. In the LLM and AI era, architects increasingly treat foundation models as production dependencies, which introduces probabilistic behavior, new security threats [14], and operational quality metrics beyond availability and latency. This evolution demands both deep technical expertise, including containerized workloads, microservices, transient infrastructure, Kubernetes orchestration, service mesh, infrastructure-as-code, observability, and practical AI and LLM systems knowledge, including retrieval architectures and AI agent architectures, and strong organizational capability for collaboration and strategic decision-making. Effective practice relies on distributed-systems patterns for consistency, fault tolerance, and resilience, and on cross-functional leadership with product, security, and business stakeholders. Skill growth is reinforced through structured design reviews, platform engineering work, incident response, open-source contributions, and internal platform development. Core technologies typically include Kubernetes, service mesh implementations, automation tooling, observability stacks, CI/CD frameworks, and AI and LLM infrastructure for model serving, vector search or retrieval components, and agent orchestration. Operating at the intersection of technology and business strategy, cloud-native architects must also define and run AI quality and safety objectives, for example factuality and hallucination reduction, robustness, privacy, and security, supported by evaluation pipelines and governance controls [11, 12, 14, 19].

## 1. Introduction

Enterprise adoption of cloud-native models has redefined the software architect from author of static blueprints to strategic leader for distributed, containerized, and transient systems. The shift from monoliths to microservices, serverless functions, and event-driven architectures demands new approaches to scalability, resilience, and delivery velocity [1, 17]. A second shift is underway: LLMs are being embedded in products and internal workflows, turning architecture into a hybrid of deterministic services and probabilistic reasoning components. This expands the architect's scope to model selection, data boundaries, and operational quality management for AI outputs. Unlike classic service dependencies, LLM failures can be silent, quality can degrade without clear error signals, and behavior can change with model updates or prompt revisions. Architects therefore need controls for evaluation and safe rollout of prompts, retrieval indexes, and model versions, comparable to how CI/CD governs code releases [12, 16].

Cloud-native platforms, such as container orchestration, service meshes, and infrastructure-as-code, introduce trade-offs, requiring deep distributed-systems expertise plus the organizational skill to drive cross-functional alignment and strategic decisions [1, 12]. Automation is foundational: IaC enables declarative provisioning, consistent deployments across environments, and fewer manual configuration errors [2]. Scalability increasingly relies on horizontal elasticity rather than vertical scaling typical of earlier architectures. Containerization improves portability and consistency, while managed containerization services automate deployment, scaling, and lifecycle management, helping systems remain available despite component failures. Effective architectures emphasize fault tolerance, observability, and rapid iteration through continuous integration and

deployment. Architects must also manage service communication, consistency tradeoffs, cost tradeoffs, and failure isolation, applying patterns such as eventual consistency, circuit breakers, and bulkheads to improve resilience under stress [1, 12]. This framework supports senior engineers transitioning into cloud-native architecture roles by summarizing the critical skills, practical development pathways, and toolsets needed to stay effective in fast-changing environments. In AI-enabled systems, automation must also include guardrails, policy enforcement, and auditability to prevent unsafe or non-compliant outputs from propagating through business workflows.

## 2. Core Technical Competencies for Cloud-Native Architecture

The basis of a successful cloud-native architecture lies in a deep comprehension of distributed systems concepts and their practical application within contemporary technology stacks. Modern architects need to showcase their ability to create systems that address the natural issues of network splits, eventual consistency, and service dependencies while taking advantage of horizontal scalability and fault isolation. Contemporary cloud-native architecture highlights essential characteristics, such as resilience, observability, and automation, that constitute the foundation of successful distributed system implementations [3]. Container orchestration is a vital skill set, with Kubernetes acting as the established benchmark for overseeing containerized tasks on a large scale. Cloud-native architectures emphasize containerization as a core design principle, allowing applications to attain portability across various environments while preserving uniform runtime characteristics [3]. Architects need to grasp not only the operational elements of cluster management but also the architectural effects of pod networking, service discovery, and resource allocation methods via thorough platform engineering techniques. Proficiency in infrastructure-as-code has become crucial for architects who need to connect development and operations teams. Cloud-native systems increasingly rely on infrastructure as code (IaC), managing and provisioning infrastructure through machine-readable definition files and reusable scripts, rather than manual configuration or interactive configuration tools [4]. Terraform, AWS CloudFormation, AWS CDK and Azure Resource Manager are tools that facilitate the codification of infrastructure specifications, enabling architects to version, test, and implement infrastructure modifications with the same level of rigor as application code via automated provisioning

workflows. The same rigor is now needed for prompts, retrieval configurations, and evaluation datasets and architects should treat them as versioned artifacts with review, testing, and controlled rollout. Service mesh technologies constitute a vital area of knowledge for handling the complexities of service-to-service communication. Contemporary cloud-native architectures utilize service mesh patterns to manage cross-cutting concerns like service discovery, load balancing, and enforcing security policies without needing changes at the application level [9]. These platforms offer essential functions for managing traffic and monitoring instrumentation while ensuring loose coupling among services via abstracted communication layers. In AI systems, a similar platform layer often emerges as a model gateway, standardizing authentication, rate limits, caching, model routing, and safety filtering, while avoiding ad hoc direct calls to model providers across teams. Metrics, logs, and traces are core observability signals in cloud-native systems, so architects should plan and instrument them early in development as part of the monitoring strategy [18]. For LLM-enabled flows, observability should capture prompt, retrieval, and tool-decision traces, and apply redaction and data minimization to reduce data leakage risk [12]. Cloud-native observability surpasses conventional monitoring methods by employing distributed tracing, real-time metrics gathering, and centralized log aggregation, offering comprehensive insight into system performance across microservices architectures. Architects need to create observability frameworks that facilitate quick incident identification and resolution, while also aiding in ongoing performance enhancement and capacity planning efforts. Scalability and elasticity are essential traits of cloud-native architectures, allowing systems to flexibly modify resource allocation in response to varying demand patterns [3, 18]. Contemporary cloud-native systems utilize horizontal scaling strategies that automatically allocate and release resources according to workload demands, guaranteeing efficient resource utilization while preserving performance levels. Architects need to create systems that enable independent scaling of individual components, facilitating detailed resource optimization and cost control within distributed application collections.

### 2.1 LLM System Design Competencies

Cloud-native architects increasingly design systems that combine microservices with LLM components such as retrieval, tool execution, and policy checks. Core competencies include retrieval-augmented

generation (ingestion, chunking, embeddings, vector search) [14], tool calling patterns (function routing, permissions, sandboxing), and fallback strategies that preserve correctness when model confidence is low. As AI agent complexity grows, developer tooling increasingly emphasizes workflow debugging and evaluation, treating agent behaviors as testable artifacts that must be regression-checked across changes in prompts, tools, retrieval, and model versions [19].

## 3. Strategic Leadership and Organizational Influence

Cloud-native architects function where technology meets business strategy, necessitating an advanced comprehension of how architectural choices influence organizational goals, team interactions, and the speed of product development. This strategic aspect reaches well past technical design, involving the skill to express the business benefits of architectural investments and synchronize technical roadmaps with changing market demands through thorough stakeholder engagement activities. In the AI era, stakeholder engagement expands to include legal, privacy, risk, and customer trust functions. Architects must define acceptable use policies, data handling rules, and escalation paths for unsafe outputs, and ensure AI features meet measurable quality and safety criteria before broad rollout. Systems thinking serves as an essential cognitive framework for cloud-native architects, allowing awareness of intricate interconnections among technology selections, organizational arrangements, and operational workflows. Decision frameworks should cover model sourcing (hosted vs self-hosted), tenant isolation for retrieval, data residency constraints, and exit strategies to reduce model vendor lock-in and revalidation costs when models change. Technical guidance and knowledge sharing are essential duties for senior architects, who need to nurture the upcoming generation of cloud-native professionals while setting architectural benchmarks and best practices throughout engineering teams. As technical professionals move up the organizational ladder, the cultivation of soft skills, such as mentoring abilities, becomes vital, with effective leadership strongly linked to the growth of interpersonal skills. This includes formulating organized methods for design evaluations, generating documentation that records technical details and reasoning behind decisions, and promoting environments of ongoing learning that allow teams to adjust to advancing technologies and practices. Collaboration skills across functions are crucial for architects, as they need to engage successfully with product managers, security teams, compliance, and business stakeholders to convert business needs into technical solutions. Soft skills such as communication and teamwork strongly affect leadership effectiveness and are often underemphasized in technical careers. This necessitates creating communication strategies that can explain intricate technical ideas to non-technical audiences while integrating business limitations and goals into architectural planning methods via systematic requirement analysis. The capability to achieve agreement on architectural choices becomes more crucial as systems become more complex and the number of stakeholders rises. Architects need to cultivate negotiation and facilitation abilities to navigate competing priorities, resource limitations, and differing technical viewpoints while ensuring sustained architectural consistency and system upkeep by utilizing thorough decision-making frameworks.

AI-driven systems raise the need for negotiation and clear communication because quality, safety, and compliance trade-offs are often ambiguous and cross-functional by nature.

## 4. Practical Pathways for Architectural Skill Development

Moving from a senior engineer to a cloud-native architect requires deliberate practice and structured exposure to progressively more complex design problems. Career-development guides and training resources suggest that candidates for cloud architecture roles must build a broad skill set that combines deep technical knowledge with business insight and leadership capabilities [5]. A useful way to structure this growth is to pair theory with recurring practice loops that resemble real architectural work, design reviews, operational learning, and platform delivery.

Design reviews are a particularly effective mechanism for developing architectural judgment. They force clarity on the true system needs, create a venue to compare competing approaches, and make trade-offs explicit across architectural styles. At the same time, strategic skill development frameworks emphasize ongoing learning in cloud-native technologies, where organizations increasingly adopt structured learning pathways that combine conceptual understanding with hands-on implementation [6]. These pathways naturally surface the kind of decisions architects repeatedly face in production, including API design, multi-tenancy considerations, and operational tooling choices.

Operational exposure is another core pathway. Incident analysis and post-mortem evaluations teach failure modes in distributed systems and reveal

whether resilience patterns work under real load and real constraints. Architects who participate in incident response and root cause analysis sharpen troubleshooting and problem-solving skills, develop intuition for system behavior under pressure, and learn to identify architectural weaknesses before they appear as production failures. Sustained improvement here depends on repeated engagement with operational challenges and disciplined documentation of lessons learned.

AI-oriented skill development fits best when it is treated as part of the same "design, operate, improve" loop rather than as a separate topic sprinkled throughout. A practical learning exercise is to build an internal RAG service end to end, including ingestion, indexing, access control, and evaluation. Security and misuse testing should be practiced explicitly through prompt-injection and tool-misuse exercises [12, 15]. Finally, the development workflow should shift toward evaluation-driven development, so that every change to a model, prompt, or retrieval component is accompanied by regression results and rollout gates, similar to how teams govern code releases [8, 10, 16]. These practices align AI work with the same reliability mindset expected in mature cloud-native systems.

Community participation and internal platform work provide additional force multipliers for skill growth. Contributing to open-source projects, especially cloud-native infrastructure, exposes architects to large-scale design constraints, collaborative development practices, and real ecosystem trade-offs [6]. Projects such as Kubernetes, Envoy, and other CNCF efforts [20] show how experienced engineers design for operability, extensibility, and safe change in complex systems. In parallel, internal platform engineering projects provide a high-leverage environment for applying cloud-native principles to real organizational needs, often requiring security controls, automation frameworks, and scalability patterns that closely reflect production realities [6]. These roles also strengthen interpersonal skills that matter for architectural leadership, because influence is exercised through alignment, mentorship, and cross-team coordination rather than direct authority.

To keep AI concerns integrated with operational readiness, architects should also run AI incident drills and treat them as first-class postmortem topics. Examples include hallucination spikes after a retrieval change, silent quality regressions after a model upgrade, or data leakage risks caused by over-permissive tool access. When these scenarios are practiced alongside standard distributed-systems incidents, the overall skill-development program becomes more coherent, and AI reliability becomes

an extension of established cloud-native engineering discipline rather than a separate track.

## 5. Essential Tools and Technologies for Modern Cloud-Native Architects

Modern cloud-native environments span a broad toolchain, so architects need depth in foundational platforms and practical fluency with the surrounding ecosystem. Contemporary cloud-native applications commonly rely on containerization, microservices, and orchestration to achieve scalability and resilience [3]. Kubernetes remains the dominant orchestration layer across clouds, serving as the de facto control plane for scheduling and operating containerized workloads. For architects, the competency is not limited to deploying workloads, it also includes understanding Kubernetes extensibility through custom resources, operators, and admission controllers, which enable policy enforcement and advanced automation at platform scale.

As microservices proliferate, service-to-service communication becomes a primary operational concern. Service mesh technologies have emerged as key infrastructure for managing these interactions, providing consistent traffic management, security, and observability at service boundaries [3]. Istio offers a broad feature set for traffic control, identity and policy enforcement, and telemetry generation, enabling standardized inter-service behaviors that can improve reliability across heterogeneous services. Lighter-weight meshes such as Linkerd are often adopted when teams want simpler operational overhead while still gaining uniform mTLS and baseline traffic and telemetry capabilities.

In some environments, organizations opt for cloud-managed service mesh offerings such as AWS App Mesh and similar provider-specific solutions. These reduce operational burden by shifting parts of lifecycle management to the cloud provider, but they can raise strategic concerns related to portability and vendor lock-in. The same portability trade-off appears in infrastructure provisioning. Infrastructure-as-code turns infrastructure changes into automated, version-controlled workflows, with Terraform commonly used as a multi-cloud declarative option for provisioning and managing infrastructure consistently. Cloud-native alternatives such as AWS CloudFormation, Azure Resource Manager, and Google Cloud Deployment Manager offer tighter platform integration, but can reduce portability and require explicit consideration in multi-cloud strategies.

Operating cloud-native systems at scale depends on strong observability and disciplined release processes. Observability platforms provide the foundation for understanding system behavior,

diagnosing failures, and identifying performance bottlenecks in distributed settings. Principles of cloud-native architecture emphasize comprehensive monitoring and logging for reliability and performance [8]. In Kubernetes ecosystems, Prometheus is widely adopted for metrics collection, often paired with Grafana for visualization and Alertmanager for routing and notification. Commercial platforms such as Datadog, New Relic, and Splunk provide unified observability suites with advanced analytics, at the cost of higher spend and potential vendor dependence.

CI/CD systems must also scale with multi-service architectures and infrastructure complexity. Cloud-native development practices rely on robust pipelines that build and deliver containerized services while automating infrastructure changes [3]. Tools such as GitLab CI/CD, GitHub Actions, and Jenkins support flexible pipeline definitions, while platforms such as Argo CD and Flux operationalize GitOps, treating Git as the source of truth for desired application and infrastructure state. GitOps aligns with cloud-native principles by emphasizing declarative configuration, versioned releases, and auditable operational control.

AI workloads introduce additional platform requirements that are best treated as a cohesive layer rather than scattered, service-specific additions. At the infrastructure level, accelerator-backed workloads often require accelerator-aware scheduling, isolation, and cost controls, implemented through practices such as separate node pools, quota management, and multi-region inference strategies. At the platform level, architects should standardize an AI layer that includes a model gateway (authentication, routing, caching, rate limiting), retrieval infrastructure for indexing pipelines and vector search [14], evaluation and prompt management (versioning, test suites) [16], release management controls (safe rollout and rollback) [8], and safety controls such as policy checks, redaction, and content filtering [10, 12]. Treating these capabilities as shared platform services reduces duplication, improves governance, and makes AI behaviors more testable and operable across the organization.

## 5.1 AI Platform Layer Tools

Modern cloud-native platforms increasingly include an AI layer composed of a model gateway (auth, routing, caching, rate limits), retrieval infrastructure (indexing pipelines and vector search) [14], evaluation and prompt management (versioning, test suites) [16], release management controls (safe rollout and rollback) [8], and safety controls (policy checks, redaction, content filters) [10, 13]. Architects should select these as platform capabilities rather than ad hoc libraries embedded inconsistently across services.

For agentic applications, orchestration frameworks such as LangGraph provide a graph-based runtime for building and operating long-running, stateful agent workflows, including durable execution and human-in-the-loop checkpoints, which helps standardize and govern tool-calling loops across teams [21].

*Table 1: Cloud-Native and LLM-Era Transition Components [1, 2, 10, 12, 13, 14]*

| Component Category | Traditional Approach | Cloud-Native Approach | Key Benefits |
|---|---|---|---|
| Infrastructure Management | Manual Configuration | Infrastructure-as-Code | Consistency, Version Control |
| Application Scaling | Vertical Scaling | Horizontal Scaling | Dynamic Resource Allocation |
| Deployment Strategy | Monolithic Releases | Containerized Microservices | Fault Isolation, Portability |
| System Resilience | Single Point of Failure | Distributed Fault Tolerance | High Availability |
| Application Behavior | Deterministic workflows | Hybrid deterministic plus LLM reasoning | Faster iteration, richer UX, new quality risks |
| Data Foundation | Relational, structured | Structured plus unstructured, embeddings, retrieval indexes | Better answers, provenance and access control required |
| Security Threats | AppSec, IAM | AppSec plus prompt injection, data exfiltration via tools | Expanded threat model, stronger governance |

*Table 2: Cloud-Native and LLM-Era Operational Attributes [3, 4, 8, 9, 10, 12, 16, 18]*

| Attribute Category | Implementation Method | Business Impact | Technical Complexity |
|---|---|---|---|

| Resilience | Fault Tolerance Patterns | High Availability | Medium |
|---|---|---|---|
| Observability | Distributed Monitoring | Operational Excellence | High |
| Automation | Infrastructure-as-Code | Reduced Manual Effort | Medium |
| Scalability | Horizontal Scaling | Performance Optimization | High |
| Evaluation | Automated offline and online evals | Quality stability, safer releases | High |
| Governance | Policies, audit logs, data boundaries | Compliance, reduced risk | High |
| Cost Controls | Token budgets, caching, routing | Predictable spend | Medium |

*Table 3: Soft Skills Impact on Technical Leadership (qualitative)*

| Career Progression Stage | Technical Skills Importance | Soft Skills Importance | Leadership impact |
|---|---|---|---|
| Entry Level | High | Medium | Medium |
| Mid-Level | High | High | High |
| Senior Leadership | Medium | Very High | Very High |
| Executive Level | Medium | Critical | Critical |

*Table 4: Cloud Architect Core Competency Areas [6,8,12,15]*

| Skill Category | Proficiency Level Required | Learning Priority |
|---|---|---|
| Technical Architecture | Expert | High |
| Business Strategy | Advanced | High |
| Leadership & Communication | Advanced | Medium |
| Security & Compliance | Expert | High |
| Automation & DevOps | Advanced | Medium |
| LLMOps & Evaluation | Advanced | High |
| AI Security & Governance | Advanced | High |

## 6. Conclusions

The evolution of software architecture roles in cloud-native contexts reflects a shift from conventional system design toward broader strategic leadership. Modern architects are expected to integrate deep technical expertise with organizational impact, and to continuously adapt to fast-changing technology environments. Success in cloud-native architecture depends on pairing technical judgment with strategic insight so architects can influence organizational direction while keeping systems sustainable, scalable, and operationally efficient over time.

On the technical side, cloud-native architects must develop advanced competency in core platforms and practices, including container orchestration, service mesh management, infrastructure automation, and comprehensive observability. They also need the architectural maturity to reason about trade-offs across reliability, security, performance, and cost, and to design systems that scale elastically while maintaining resilience under operational stress. The expanding cloud-native tool landscape reinforces the need for both depth in foundational platforms and breadth across supporting technologies such as Kubernetes, service mesh solutions, infrastructure-as-code tools, observability platforms, and CI/CD systems.

In parallel, the role demands leadership capabilities that extend beyond implementation choices. Effective architects drive cross-functional collaboration, align technical decisions to business goals, shape platform direction, and establish shared standards that teams can implement consistently. In the LLM era, this leadership mandate expands further: evaluation and governance become first-class architectural pillars alongside scalability and observability, because quality and safety failures can be as damaging as traditional outages [11, 13]. Architects therefore need to institutionalize practices that keep model behavior measurable, changes controlled, and risks managed as part of normal delivery.

Practical skill development is best approached through repeated exposure to real architectural work. Structured growth paths that emphasize participation in design reviews, platform engineering initiatives, incident response and postmortems, and open-source engagement create the experience base required for

mature architectural judgment. These pathways strengthen both the technical competencies needed to build and operate distributed systems and the interpersonal capabilities required to lead effectively in complex organizations.

## Author Statements:

## References

[1] Lee Atchison, "What You Need to Learn to Become a Cloud-Native Architect," Cloud Native Now, 1 November 2022.

[2] Hemanthnvd, "Embracing the Cloud-Native Mindset: Delving Deep into Infrastructure, Automation, and Scalability with CSYE 6225," Medium, 24 April 2024. Available: https://medium.com/@hemanthnvd/embracing-the-cloud-native-mindset-delving-deep-into-infrastructure-automation-and-scalability-f4c489da1c1f

[3] Azure Architecture Center, "Cloud design patterns" (updated July 18, 2025). Available: https://techcommunity.microsoft.com/blog/appsonazureblog/step-by-step-practical-guide-to-architecting-cloud-native-applications/4057960

[4] Indika Kumara, et al., "The Do's and Don'ts of Infrastructure Code: a Systematic Grey Literature Review," Information and Software Technology, September 2021. Available: https://www.sciencedirect.com/science/article/pii/S0950584921000720

[5] Sneha Chugh, "How to Become a Cloud Architect: Top 10 Skills to Master," Emeritus, 14 November 2024. Available: https://emeritus.org/blog/how-to-become-cloud-architect/

[6] CertLibrary, "Developing a Comprehensive Strategy for Cloud Native Skills Growth,". Available: https://www.certlibrary.com/blog/developing-a-comprehensive-strategy-for-cloud-native-skills-growth/

[7] NIST Special Publication 800-233 (2024), "Service Mesh Proxy Models for Cloud-Native Applications." Available: https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-233.pdf

[8] Tom Grey, "5 principles for cloud-native architecture - what it is and how to master it," Google Cloud, 20 June 2019. Available: https://cloud.google.com/blog/products/application-development/5-principles-for-cloud-native-architecture-what-it-is-and-how-to-master-it

[9] National Institute of Standards and Technology, "Artificial Intelligence Risk Management Framework (AI RMF 1.0)," NIST, 26 January 2023. Available: https://doi.org/10.6028/NIST.AI.100-1

[10] C. Autio, R. Schwartz, J. Nadeau, K. Grama, A. Hsiang, H. Nguyen, and K. Roberts, "Artificial Intelligence Risk Management Framework: Generative Artificial Intelligence Profile," NIST, July 2024. Available: https://doi.org/10.6028/NIST.AI.600-1

[11] ISO/IEC, "Information technology, Artificial intelligence, Guidance on risk management (ISO/IEC 23894:2023)," ISO/IEC, 6 February 2023. Available: https://webstore.iec.ch/en/publication/82914

[12] OWASP Foundation, "OWASP Top 10 for Large Language Model Applications, 2025," OWASP, 18 November 2024. Available: https://owasp.org/www-project-top-10-for-large-language-model-applications/assets/PDF/OWASP-Top-10-for-LLMs-v2025.pdf

[13] Werner Vogels (Amazon CTO), "Return of The Frugal Architect(s)" (Dec 5, 2024, All Things Distributed) Available: https://www.allthingsdistributed.com/2024/11/return-of-the-frugal-architect.html

[14] Patrick Lewis, Ethan Perez, Aleksandra Piktus, et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," Advances in Neural Information Processing Systems (NeurIPS 2020), 2020. Available: https://papers.nips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html

[15] Jun Yan, Vikas Yadav, Shiyang Li, et al., "Backdooring Instruction-Tuned Large Language Models with Virtual Prompt Injection," Proceedings of NAACL-HLT 2024 (Volume 1: Long Papers), June 2024. Available: https://aclanthology.org/2024.naacl-long.337/

[16] Percy Liang, Rishi Bommasani, Tony Lee, et al., "Holistic Evaluation of Language Models," Transactions on Machine Learning Research, 2023. Available: https://openreview.net/forum?id=iO4LZibEqW

[17] CNCF Authors, "Level 1, Build," Cloud Native Maturity Model, last modified Sep. 25, 2025. [Online]. Available: https://maturitymodel.cncf.io/level-1/

[18] Google Cloud, "Ensure operational readiness and performance using CloudOps." https://docs.cloud.google.com/architecture/framew ork/operational-excellence/operational-readiness-and-performance-using-cloudops

[19] Victor Dibia, et al., "AutoGen Studio: A No-Code Developer Tool for Building and Debugging Multi-Agent Systems," arXiv, 2024. Available: https://arxiv.org/abs/2408.15247

[20] Cloud Native Landscape, "Application Definition & Image Build," Available: https://landscape.cncf.io/

[21] Langchain Docs, "LangGraph overview." https://docs.langchain.com/oss/python/langgraph/o verview