

# SERIES CHEAT SHEET

*Production is not deployment. It is the architecture of trust under load.*

AUTHOR: NABEEL A. KHAN · THREE BOOKS, THREE REFERENCE ARCHITECTURES, NINE RESEARCH ANCHORS · COMING SEPTEMBER 2026 · NABEELKHAN.COM/SERIES-RESOURCES

Three books, one company, one discipline observed from three altitudes. Set inside **Nebula Financial**, a regulated fintech running three systems, each owning a layer of one stack. A single fraud signal traces through all three and proves why each layer must exist: NexusCore routes it, AgentMesh investigates it, ThinkFlow ships the fix.

## THE THREE BOOKS AT A GLANCE

BOOK	SYSTEM	LAYER	PRIMARY READER	WHAT IT TEACHES
<b>Book 1</b> <i>LLM Systems in Production</i>	NexusCore	Infrastructure / SRE	SREs, cloud architects, platform engineers entering AI	How a request finds the right model under a latency, cost, and risk budget, and is logged as evidence.
<b>Book 2</b> <i>Prompt Systems &amp; Agent Orchestration</i>	AgentMesh	Application / Agents	AI engineers, application developers, agent builders	How a pile of prompts becomes governed agents that plan, act, verify, and submit to human review.
<b>Book 3</b> <i>DevOps for AI-Native Platforms</i>	ThinkFlow	Operations / Platform	DevOps and platform teams, internal-platform owners	How a platform builds, ships, and governs the models, agents, and code the other two layers depend on.

## THE THREE REFERENCE ARCHITECTURES IN BRIEF

<p>NEXUSCORE · BOOK 1</p> <h3>Routing and observability gateway</h3> <ul style="list-style-type: none"> <li>■ <b>Data.</b> Gateway to Routing Brain to Latency Controller to the governed model pool (on-prem, cloud, frontier).</li> <li>■ <b>Control.</b> The Routing Brain chooses the smallest model a request can trust; the Latency Controller earns its SLO with speculative decoding.</li> <li>■ <b>Governance.</b> Routing and speculation policies as signed, versioned artifacts promoted through GitOps with a conformance gate.</li> <li>■ <b>Evidence.</b> Every routing decision, decode strategy, model version, and cost recorded as a queryable record.</li> </ul> <p><i>The gateway is where an institution decides what it is allowed to think, and proves it thought it lawfully.</i></p>	<p>AGENTMESH · BOOK 2</p> <h3>Catalog of record and orchestration</h3> <ul style="list-style-type: none"> <li>■ <b>Taxonomy and contracts.</b> Agents classified by domain, capability, and autonomy, each declaring a capability contract.</li> <li>■ <b>Four-module loop.</b> Planner decomposes, executor calls tools under contract, verifier checks output, generator responds.</li> <li>■ <b>Graph orchestration.</b> Workflows as activity-on-vertex graphs that can be inspected, parallelized, and refined.</li> <li>■ <b>Human review as runtime.</b> Tiered review (auto, lightweight, full supervisory) mapped to risk; feedback is first-class data.</li> </ul> <p><i>An agent is not a clever prompt. It is a governed actor with a contract, a boundary, and a record.</i></p>	<p>THINKFLOW · BOOK 3</p> <h3>AI-native developer platform</h3> <ul style="list-style-type: none"> <li>■ <b>Catalog and golden paths.</b> Services, models, and agents as first-class catalog citizens, with paved roads for AI workloads.</li> <li>■ <b>Policy-bounded delivery.</b> Agentic gates inside the pipeline, constrained by Delivery Guardrails and a trust-tier model; every decision logged.</li> <li>■ <b>RL-learned testing.</b> An adaptive agent chooses which tests to run, skip, or parallelize, bounded by per-service risk profiles.</li> <li>■ <b>PARA operations.</b> DevOps agents as perception, action, reasoning, reflection, each evaluated against benchmarks of real failures.</li> </ul> <p><i>A platform is the paved road that decides what an organization can build without losing control.</i></p>
--	---	---

## THE NINE RESEARCH ANCHORS

#	BOOK	ANCHOR	CH.
1	NexusCore	Universal, workload-aware LLM routing	B1, 4
2	NexusCore	SLO-aware speculative and pipelined decoding	B1, 6
3	NexusCore	Secure and auditable router lifecycle	B1, 9
4	AgentMesh	Trainable planner-executor-verifier-generator loops	B2, 3
5	AgentMesh	Graph-based agent workflows with dynamic refinement	B2, 5
6	AgentMesh	Tiered human-in-the-loop orchestration	B2, 8
7	ThinkFlow	Policy-bounded AI-augmented CI/CD with trust tiers	B3, 4
8	ThinkFlow	Reinforcement-learned adaptive test selection	B3, 6
9	ThinkFlow	Perception-action-reasoning-reflection DevOps agents	B3, 7&8

## THE SPINE IN FIVE CHAPTERS

- 1 Book 1, Ch. 4.** The router selects a model worthy of the risk and logs the decision.
- 2 Book 1, Ch. 12.** The routed request is handed forward as a workflow.
- 3 Book 2, Ch. 3.** The workflow is decomposed across planner, executor, verifier, generator.
- 4 Book 2, Ch. 8 & 11.** Tiered human review engages and the workflow produces an audit record.
- 5 Book 3, Ch. 11.** The gap the investigation exposed becomes a new service the platform ships.

## THE GOVERNANCE SPINE

TOGAF for architecture governance. DMBOK for data governance. Routing policies become governed artifacts with provenance and promotion. Agent catalogs become managed assets with lineage. Delivery decisions become audited events. Every book's Appendix E carries the TOGAF, DMBOK, and EU AI Act / NIST AI RMF alignment checklists.

*Governance is not compliance. It is coherence made visible.*



FROM THE PAGE TO A SYSTEM

### The template package.

Turn the series artifacts into governance-grounded implementation starting points, adapted to your stack and compliance obligations. Scan the code or visit [nabeelkhan.com/engagements](https://nabeelkhan.com/engagements) to start an engagement. The full guides and print-ready downloads live at [nabeelkhan.com/series-resources](https://nabeelkhan.com/series-resources).