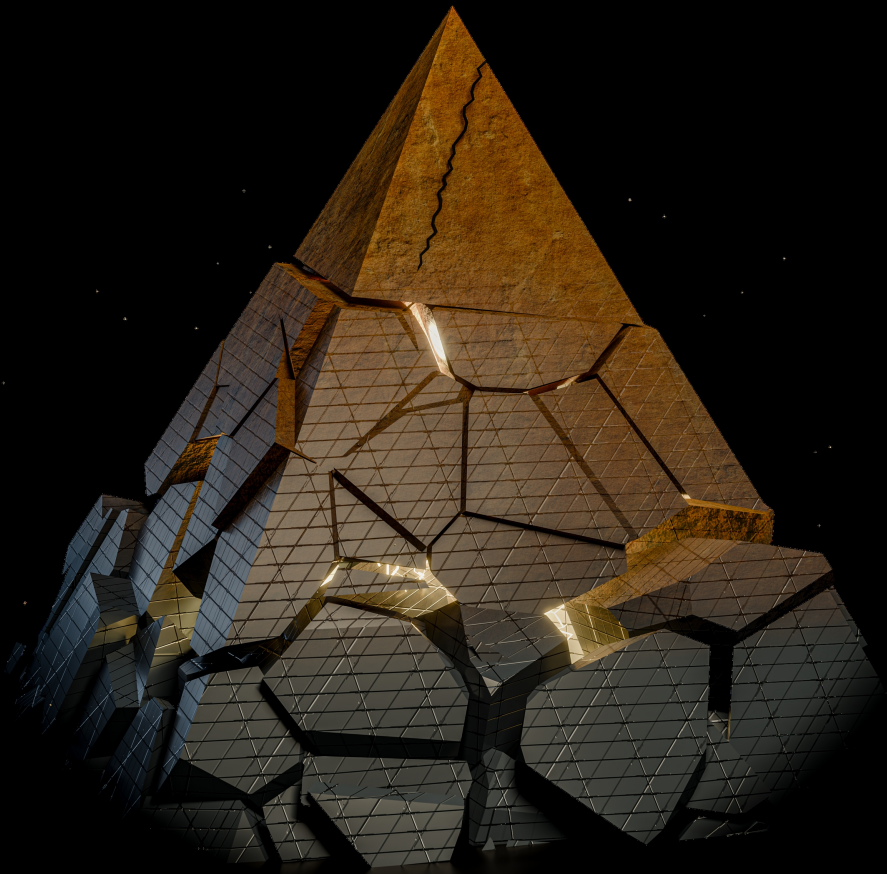


BUILDING AI-NATIVE PROFESSIONAL SERVICES FIRMS

Strategy, Economics, and Execution



DANIEL M. KATZ • MICHAEL J. BOMMARITO II • JILLIAN BOMMARITO

CHAPTER 5

Thin Wrappers and Deep Systems

CHAPTER 5 · APRIL 2026

- ✓ Prologue · Four Lawyers, One Monday Morning

PART I · THE OPPORTUNITY

- ✓ Chapter 1 · The AI-Native Professional Services Firm
- ✓ Chapter 2 · The Economics of Transformation
- ✓ Chapter 3 · Uncertainty and Structure as Strategy

PART II · THE TRANSFORMATION

- ✓ Chapter 4 · Craft, Process, and Scale
- **Chapter 5 · Thin Wrappers and Deep Systems** YOU ARE HERE
- Chapter 6 · Land, Expand, Compound
- Chapter 7 · The Org Chart No One Expected
- Chapter 8 · Growing Pains: Talent and Culture
- Chapter 9 · Measure Twice, Cut Once

PART III · THE EXECUTION

- Chapter 10 · Paths—Build, Transform, Combine
- Chapter 11 · Capital and Investment
- Chapter 12 · Execution

PART IV · THE VALUE

- Chapter 13 · Value Creation and Capture
- Chapter 14 · The Future

Chapter 5: Thin Wrappers and Deep Systems

We shape our tools, and thereafter our tools shape us.

— Marshall McLuhan

5.1 FUNDED BUT FRAGILE

The seed round had finally closed on a Tuesday in the late summer—several months after Sarah had shaken hands with Alex and the other investors in various Bay Area coffee shops and nondescript buildings. Several months for three million dollars in seed capital. The lawyers on both sides had treated it like the InBev–Anheuser-Busch merger, negotiating protective provisions and investor rights as though the fate of a Fortune 500 company hung in the balance rather than the operating account of a small law firm. Sarah had bitten her tongue through most of it, but she told David one morning: “This is why our clients hate lawyers. My own seed round needed more legal documentation than half the deals I closed at my old firm.”

The wire hit in two tranches, the first landing in the firm’s operating account at 9:47 AM Pacific. Alex Hawthorne’s \$500,000 anchor had done exactly what he promised it would: it attracted two additional angel investors for another \$500,000, and a seed fund that Maya connected her to filled out the remaining \$2 million. Three million dollars at a \$15 million post-money valuation. Sarah

watched the balance on her screen for longer than she would have admitted to anyone. Not all that long ago she had been a senior associate billing more than two thousand hours a year at a firm that treated AI like a parlor trick. Now she had capital, a small team, a regulatory structure that let her do things most traditional firms could not or would not, and—most pressingly—a technology problem that no amount of operational discipline alone could solve.

David Park had been on board since earlier in the summer and was already reshaping how the firm operated: building the production system, the process maps, the quality protocols that were transforming how the firm delivered work. He had not waited for the money to arrive. Neither had Sarah. The work he had done in those months before the wire hit proved that operational discipline did not require capital, just clarity. His process maps covered three whiteboards. His cycle-time measurements had exposed eleven handoff points where work sat idle. Run-of-the-mill contract review turnaround had dropped from five days to two, and batch analysis was markedly quicker. Error rates fell. Client satisfaction ticked upward.

But David's improvements had also exposed a harder truth: the workflows he was improving ran on duct tape. The AI analysis was solid—genuinely solid—but the infrastructure connecting everything was fragile, manual, and held together by Google Sheets, Notion pages, and the collective memory of a small team of people. David could refine the processes. He could not engineer the platform. For that, Sarah needed someone who understood both legal technology and software architecture at production scale. She needed a proper CTO.

Before the round closed, the search had gone nowhere. The candidates Sarah could attract were freelancers and consultants, competent people but not the kind of technical leader who could architect a proprietary platform from the ground up. A law firm at this size, however ambitious its thesis, was not where serious engineers went to build careers. But a funded startup with three million in the bank and a fifteen-million-dollar valuation—that was a different conversation entirely. The capital did not just buy runway. It bought credibility.

It told candidates that real investors had evaluated this opportunity and written real checks. It meant their generous equity grants might actually be worth

something. It meant the resources to build properly rather than perpetually patching. Not just thin wrappers but a chance to build deep systems.

The search still took longer than Sarah wanted. Through August, she had interviewed a dozen candidates. Some were gifted engineers who had never touched a legal document and treated the domain as a straightforward NLP problem. “Just tokenize the contracts and let the model sort it out,” said one candidate. Others were legal technology veterans whose ambition extended no further than building a marginally better document management system.

She needed someone who lived in the gap between the two worlds, someone who understood that legal documents were not just text to be parsed but instruments of consequence. A contract allocated risk between counterparties. A regulatory filing calibrated candor against exposure. A motion reshaped the facts a judge would see. Each carried economic weight and professional liability.

Alex made the introduction in early September. Priya Chandrasekaran had spent three years at Ledgerwell, a fintech startup that built compliance and risk-scoring tools for financial institutions. The work sat close enough to law that Priya had learned to think in terms of regulatory frameworks, contractual obligations, and audit trails. She had risen from a senior developer to lead engineer on their transaction monitoring platform, and she had watched that company build features before architecture, bolt on integrations after the fact, and accumulate technical debt that eventually slowed development to a crawl. She was ready to build something from the ground up, and she understood, from hard experience, what that actually required.

Sarah knew within twenty minutes of their first conversation. Priya asked questions nobody else had asked: What is your data retention policy for client documents? How do you handle model versioning when a frontier lab releases an update? What happens to your workflow if the API goes down for four hours? How much preprocessing are you doing before legal documents hit the model—are you extracting structure from the raw documents or just feeding in flat text?

That last question revealed something the others had missed. Her years building compliance tools at Ledgerwell had taught her that the intelligence of any system was bounded by the quality of its inputs, and legal documents, with

their nested definitions, cross-references, and exhibit structures, demanded more care than most. The other candidates had talked about what they could build. Priya talked about what could break.

Sarah offered her the title of CTO, grand for what was still an eight-person company but a signal of something Sarah increasingly saw as non-negotiable: that technology leadership would be co-equal with legal leadership, not subordinate to it. This was a different version of Sarah. Not all that long ago she had succumbed to the cult of law and lawyers, had seen them as the center of the universe. She was not that person now. She had a more expansive view, and she was ready to meet the moment.

Priya accepted Sarah's offer, but she could not start immediately. She had a notice period at Ledgerwell and two projects she refused to leave unfinished, a professional discipline that only made Sarah want her more. They agreed on an early October start date. Sarah spent the intervening weeks watching the duct-tape infrastructure strain under growing client volume, each passing day reinforcing the urgency of what Priya would find when she arrived.



5.2 A TANGLED WEB OF DUCT TAPE

Sarah stood behind Priya at the whiteboard, watching her diagram the firm's technical architecture with the methodical precision of someone disassembling a crime scene. It had been three weeks since her start date, and Priya had spent the first two conducting what she called a "full-stack audit": mapping every API call, every prompt template, every spreadsheet, every Slack message that constituted the firm's technology infrastructure. Now she was presenting her findings, and Sarah could tell from the set of Priya's jaw that the news was not going to be flattering.

David Park sat at the conference table, his black coffee untouched beside his signature Moleskine notebook and Leuchtturm 1917 pen he carried everywhere. As COO, he had spent the past months building the operational discipline that

turned Sarah’s vision into repeatable processes. Now he was waiting to hear whether the technology could keep up.

“Okay,” Priya said, capping her marker. “Here is what we have.”

The whiteboard showed a tangle of arrows connecting boxes with labels like *Claude API*, *Google Sheets QA Tracker*, *Prompt Library (Notion)*, *Manual Email to Client*, and *Airtable Intake Form*. Red lines connected components that had no formal integration: data that moved between systems via copy-paste or, in one case, a screenshot someone texted in a group chat.

“What you have built,” Priya said, “is the technological equivalent of an airplane held together with duct tape and good intentions. It works. I want to be clear about that—it works. You are delivering real value to clients, and the core AI analysis is genuinely good. But the infrastructure connecting everything is fragile, manual, and completely unscalable.”

David leaned forward. “How fragile?”

“If someone changes the format of a prompt template in Notion, it breaks the downstream extraction pipeline. If the Claude API has a rate limit issue, there is no retry logic; someone has to manually resubmit. Your QA process lives in a Google Sheet that three people edit simultaneously, and last week two reviewers overwrote each other’s corrections on one recent engagement.” Priya turned to Sarah. “You told me you wanted to 10x the firm’s capacity within eighteen months. On this infrastructure, you will be lucky to 2x without things falling apart.”

Sarah absorbed this without flinching. She had known it was bad. She had not known it was *that* bad. When you are building a plane while flying it, you do not always notice how many parts are held on with zip ties.

“So what do we do?” Sarah asked.

Priya uncapped a fresh marker. “First, we figure out what to build and what to buy. That is the decision that determines everything else.”

 **BUILD OR BUY**

Buy: Commodity capabilities that everyone needs but nobody differentiates on: foundation models, practice management, billing, security tools. Some purchases are off-the-shelf subscriptions; others are deeper integrations with specialized vendors for capabilities like document processing or OCR. Either way, the logic is the same: if it is not core to differentiation, do not build it.

Build: Capabilities that sit between the foundation models and the commodity tools: the retrieval layer, workflow orchestration, model harness, quality systems, and analytics. This is where the moat is.

Every AI-native firm faces the same question: which technology capabilities to develop internally and which to purchase from vendors. The answer, Priya argued, depends on where competitive differentiation lies. Buy commodity infrastructure that everyone needs but nobody differentiates on, whether that means a straightforward subscription or a deeper integration with a specialized vendor. Build the capabilities that create sustainable advantage: the layers that embody the firm’s institutional knowledge and operational discipline, the things that would make Candor’s work materially better than a competitor using the same underlying models. “These boundaries are not permanent,” Priya added. “A capability worth buying today may be worth building in a year, once we understand our requirements and have the engineering capacity. We can revisit this quarterly.”



The build-versus-buy conversation consumed most of the next week. Sarah had been operating on a simple heuristic since launch: use the best available tools, connect them with whatever works, and focus on delivering quality to clients. That approach had gotten the firm from zero to eight people and a growing roster of mid-market corporate clients. It was not going to get them to fifty.

Priya started by cataloging every technology capability the firm needed, then sorting them into two buckets. She drew a table on the whiteboard during their Monday morning meeting.

“Foundation models, Claude, GPT, Gemini, whatever we decide to use for specific tasks, we buy. There is no universe in which we train our own large language model. The cost would be absurd and the result would be inferior. We are not OpenAI, Anthropic, or Google. We are not trying to be a frontier lab. We use their models the way a construction company uses steel—we buy the raw material and build with it.”

David nodded. “That is commodity infrastructure. We need it, everyone needs it, and there is no competitive advantage in producing it ourselves.”

“Exactly.” Priya moved to the next column. “Document storage, billing, security infrastructure, we buy those too. There are mature vendor solutions. Building our own billing system would be a waste of engineering time.”

“And the second bucket?” Sarah asked, though she suspected she knew the answer.

“Everything in between.” Priya circled a large section of the whiteboard. “The retrieval layer, how we find and surface relevant information from client documents and our own knowledge base. The workflow orchestration, how work moves through AI and human stages. The application layer, the interfaces our lawyers and clients interact with. The quality system, how we monitor and measure and improve AI output. And the analytics layer—how we learn from every engagement.” She set down the marker. “That is where our moat is. That is what we build.”

Sarah sat with this for a moment. The logic was clean, but the implications were expensive. Building all of that required engineers, infrastructure, time, and capital. She had spent eight years at a law firm where technology meant a new version of Microsoft Office, and that instinct told her there had to be a shortcut.

“Before we commit to building all of that,” Sarah said, “I want to test what is already on the market. There are like a million vertical legal tech products designed for exactly what we do. If one of them is good enough, we save ourselves a year of development.”

Priya's expression was carefully neutral. "That is a reasonable thing to test."

David, ever the process person, suggested a structured evaluation. "Short pilot. Clear criteria. Documented results. Then we decide."

5.3 THE WRAPPER AND THE MOAT

ContractZoo.ai was one of the most talked-about vertical legal AI platforms in the market. Sarah had seen their name in every legal tech newsletter for the past six months: polished interface, well-funded team, a client list that included several Am Law 200 firms as well as Fortune 500 legal departments. When she reached out to schedule a demo, the founder himself got on the call.

James Martin joined the video from what appeared to be a converted garage in Palo Alto—the kind of space that signaled startup authenticity in Silicon Valley. He was younger than Sarah expected, probably late twenties, with the restless energy of someone who had recently discovered that venture capital could turn an idea into a company.

James's path to ContractZoo had been anything but direct. He had practiced law for a couple of years at a large law firm in San Francisco, handling corporate matters and hating almost every minute of it—not the substance, which he found intellectually engaging, but the practice itself. The timesheets. The partner politics. The soul-crushing realization that his days consisted of marking up the same boilerplate provisions in slightly different documents for slightly different clients.

He quit to apply to Y Combinator with an idea for a matching platform that would connect clients with lawyers the way Airbnb connected travelers with hosts: profiles, ratings, algorithmic matching based on practice area and budget. Despite having far from an original idea, YC accepted him. The matching platform did not survive its first contact with reality. Lawyers did not want to be rated like Uber drivers. Both enterprise and retail clients were dubious about trusting an algorithm to pick their lawyer. The unit economics required demand that simply did not exist.

His roommate and technical co-founder, Kevin, convinced him to pivot. Their next idea was LegalPulse, a platform that used natural language processing

to analyze court filings and predict case outcomes—a kind of Moneyball for litigation. The pitch was compelling in the abstract but data access proved challenging and their ability to leverage NLP/ML models available at the time was insufficient to deliver on the idea. The predictions were marginally better than a coin flip, and no lawyer was going to stake a client’s case on marginally-better-than-random odds. They burned through most of their YC timeline before James pulled the plug on that idea.

The third pivot was the one that stuck, and it stuck largely because of timing. Around late 2022 and into early 2023, the frontier language models, Claude and GPT-4, had crossed a threshold that the previous generation had not. They could actually read a contract, identify the relevant provisions, and produce analysis that was coherent and largely accurate. James, who understood contracts from his practice years and understood the market from his failed startup attempts, recognized the opportunity immediately.



He and Kevin built ContractZoo in three months: a clean interface wrapped around the new models, with legal-specific prompt templates that structured the analysis. The initial product was built in Replit and it was pretty good. More honestly, the product was good because the underlying models were good, and the models had become good at exactly the right moment for James to build a business on top of them. James was riding the AI wave hoping it would not crest.

Other than the part about being a thin wrapper, he recited most of his back story within the first few minutes of the call with Sarah and team. He narrated his journey with the practiced cadence of someone who had delivered the story to dozens of investors and customers. He framed the failed pivots as “learning experiences” and the timing as “strategic patience.” Sarah, who had her own complicated relationship with the gap between startup mythology and operational reality, listened with a mix of recognition and wariness.

“We built ContractZoo because I saw firsthand how much time lawyers waste on contract review that should be automated,” James said. He shared his screen and walked Sarah, Priya, and David through the platform’s interface. “Our *proprietary legal intelligence layer*TM handles the analysis. The lawyer uploads the documents and gets structured output. That is the whole workflow.”

The interface was genuinely impressive. Clean design, intuitive navigation, real-time processing indicators. James fed a sample NDA through the system and the analysis appeared in under a minute: parties identified, key terms extracted, risk flags color-coded by severity. The output looked professional. It looked like something a client would trust.

“What models are you running underneath?” Priya asked. She had her notebook open but had not written anything yet.

“We use a combination of frontier models,” James said. “We have what we like to call a proprietary blend of frontier models. But on top of this is our legal layers. We have built *legal-quality*TM intelligence layer that goes beyond what you would get from a raw model.”

“Can you describe the analysis layer?” Sarah asked. Her computer science background had taught her to listen for the gap between marketing language and technical reality. *Proprietary legal intelligence layer*TM was a phrase that could mean almost anything. *Lots of unenforceable trademarks being asserted—kinda feels like a smokescreen.* Sarah thought this to herself but politely did not say. James was selling hard and right below the surface she felt a degree of desperation that was not so far from view.

“Sure. We have developed custom information architectures and prompting taxonomy as well as a series of preprocessing techniques. We have an agentic layer which orchestrates a series of our technologies including our specialized extraction module which flags risks specific to that contract category. We also have a post-processing layer that normalizes the output into consistent formats.”

Sarah had now counted three layers. For a zoo, she thought, it had an awful lot of gift wrapping and not much wildlife. She waited a beat. “So the analysis layer is prompt engineering and output formatting.”

James's smile held, but something shifted behind it. "It is significantly more sophisticated than that. We have spent two years refining the prompts, building the templates, training the post-processing—"

"I do not doubt the effort," Sarah said. She kept her tone neutral—she was not trying to embarrass him. "But I want to understand the architecture. If I took the same contract, wrote a detailed prompt specifying the analysis framework, and sent it directly to one or more frontier models, would the output be materially different from what your platform produces?"

The pause was brief but telling. "The output would be less consistent," James said. "Our layered approach ensures consistency across thousands of contracts. That is the value: you do not have to write the prompts yourself, manage the API calls, format the output. We handle all of that."

Priya made a note in her notebook. Sarah could see the single word she wrote: *wrapper*.

"One more question," David said. He had been listening with his arms folded, the posture he adopted when he was in evaluation mode. "When a client catches an error in the AI's analysis and corrects it, does that correction feed back into the system? Does the platform learn from it?"

"We collect feedback through our support process," James said. "Our team reviews it and incorporates improvements into future template updates."

"How often do the templates update?"

"Quarterly." James paused, perhaps sensing the direction of the question. "We are working on more frequent update cycles. One thing to note is that some clients are hesitant to allow for this specific feedback loop to be constructed because of client confidentiality. So as of right now we are not able to support detailed learning across customers or even across the clients of a given customer. However, we do take general feedback from users and incorporate it into our products."

David looked at Sarah. This was not a minor point. The data from every client's corrections were not looping back into making ContractZoo better. Instead, there was only a tangential flywheel of general product improvements. The flywheel David had been designing—where every engagement made the AI

better for the next one—did not exist in ContractZoo’s architecture because the clients would not allow for their product vendor to hold their legal work product.

Sarah thanked James for the demo. He was gracious, offered a free pilot period, and followed up with pricing that was genuinely competitive: \$95,000 per year for unlimited contract analysis. “We have a \$20k one month starter package,” said James. “You do not have a free trial,” asked Sarah. “No, we are just so busy that we have a minimum just to get started. However, we would credit you the \$20k if you sign for one year thereafter.” She told him she would talk to her team and get back to him within a few weeks.

After the call, the three of them sat in the conference room for a moment without speaking. Priya broke the silence.

“Listen, he built a front end on top of someone else’s intelligence.”

“To be fair,” Sarah said, “the front end is genuinely good. And the pricing is aggressive. For an in-house legal team that does not have its own AI workflow, ContractZoo saves them real time and effort.”

“For an in-house team, yes,” Priya said. “For us, it is a step backward. We would be handing our workflow to someone else and getting nothing back. Building someone else’s platform and not our own.”

David nodded. “But we should prove that with data, not intuition. Run the pilot. Let the numbers make the case.” Sarah told Priya she generally agreed but she saw the question slightly differently. “Maybe we use ContractZoo for a year while we build our own system and begin to transition over sometime in the next 12 months. But let’s just see,” said Sarah.

Sarah ran the pilot herself. She selected three recent engagements where the firm had already delivered work to clients, so she could compare ContractZoo’s output against their own. She fed the same documents through the platform, used their templates, followed their recommended workflow. She also trialed two other platforms, one focused on due diligence and another on regulatory compliance review, to see if the pattern held.

By day four, she had her first concern. The outputs were competent but generic. ContractZoo’s contract analysis identified the standard risk categories

(indemnification, limitation of liability, termination provisions, intellectual property assignment) with reasonable accuracy. But the analysis lacked the specificity that her clients cared about. When Sarah's team reviewed a vendor agreement for a mid-market SaaS company, they flagged a subtle interaction between the auto-renewal clause and the price escalation provision that created a compounding cost exposure over a five-year term. ContractZoo's analysis noted both provisions independently but missed the interaction entirely.

By day eight, the pattern was clear. Sarah wrote up her assessment in the shared document David had created for the pilot.

ContractZoo produces output that is roughly equivalent to what a competent associate would produce on a first pass, which is to say, roughly equivalent or perhaps even marginally better than what GPT or Claude produces with a well-designed prompt. The UI is polished. The workflow is clean. But the underlying analysis is a thin wrapper around the same frontier models we already use, with rigid templates that cannot handle the edge cases our clients actually care about.

The platform has three specific limitations that are disqualifying for our use case. First, the templates are fixed. We cannot modify the analysis framework for different client contexts. A SaaS vendor agreement requires different risk weighting than a manufacturing supply contract, and ContractZoo treats them identically. Second, the system has no mechanism for learning from our corrections. When we identify something the AI missed, that knowledge stays in our heads. It does not feed back into the platform. Third, the client-facing output is formatted for their templates, not ours. Our clients expect deliverables that reflect our firm's analytical framework, not a generic vendor report.

She shared the assessment with Priya and David at their Wednesday standup.

⚠ THE THIN WRAPPER TRAP

Products that wrap frontier models with prompt templates and a polished UI, without proprietary data, learning loops, or workflow customization, are vulnerable to commoditization. If a well-funded competitor can replicate the offering in months, there is no moat. The test: would sending the same document directly to the underlying model with a well-designed prompt produce materially different output? If the answer is no, the product is a thin wrapper.

“We are not moving forward after the pilot but we might as well just finish the pilot period,” Sarah said. “The product is not bad. It is actually pretty decent. But it is not better than what I think we can build ourselves, and this will constrain us in ways that will matter more as we grow.”

Priya, who had been reading the assessment on her laptop, looked up. “I want to make sure the reasoning is right, because this is the decision that determines the next years of my life.” She paused. “The products you tested are wrappers. They take the same foundation models we use, add some legal-specific prompt templates, and put a nice interface on top. The value is convenience: you do not have to build the workflow yourself.”

“Right.”

“But convenience is not defensibility. If ContractZoo’s competitive advantage is their UI and their templates, any reasonable competitor—including us—could theoretically replicate that in months. There is no moat. The moat has to come from somewhere else.” Priya stood up and walked to the whiteboard again. “Buy the foundation models. Build the workflows, the knowledge layer, and the quality system. That is where our moat is.”

David had been quiet through the exchange, but now he spoke. “I want to add something from the operations side. The biggest problem with ContractZoo is not the AI quality. It is the data. Every contract we review through their platform generates learning, corrections, edge cases, client preferences, and that learning stays in their system, not ours. If we build our own pipeline and can

get the right sort of client consent to data aggregation, that data is our asset. It compounds. It makes us better over time in ways that a vendor platform never will.”

Sarah looked at both of them. The CTO saying build. The COO saying build. The pilot confirming that buying was not good enough. The decision was not hard. But executing upon that decision would be.

“Okay,” she said. “We build.”

The decision applied beyond just contracts and ContractZoo. Over the subsequent months, Sarah would end up evaluating half a dozen specialized legal AI point solutions, and the pattern was consistent. Each offered a polished interface wrapped around the same foundation models her firm already used. Each marketed proprietary intelligence that, under scrutiny, amounted to little more than prompt templates and output formatting. Each charged subscription fees for convenience that a well-configured general-purpose model could roughly match. And none could keep pace with the frontier: every time the foundation model providers released a new version, the gap between what the raw model could do and what the wrapper added narrowed further. Building on a wrapper meant waiting for someone else’s engineering team to decide which capabilities mattered. Building on the foundation meant every advance was available on day one and could be combined directly with the services offered by Candor without any additional vendor markup.

Sarah emailed James Martin a couple of weeks later. She kept it brief and honest: the platform was well-built, but her firm needed the ability to customize workflows and capture correction data in ways that a third-party product could not accommodate. She wished him well.

James’s reply came an hour later: “*Understood. If your needs change, we are here.*” Two sentences. Professional. But Sarah wondered, reading them, whether he heard the subtext as clearly as she had intended it. His product was good. It was just not different enough from what the underlying models already provided. And if she could see that after a two-week pilot, his customers (and investors) would eventually see it too.

What Sarah did not know—could not have known—was that James had been asking himself the same question for months. He had spent the better part

of the past year on the interface: the clean design, the intuitive navigation, the real-time processing indicators, the color-coded risk flags that made the output feel authoritative. The UI was genuinely excellent. He was proud of it. And from a sales and marketing perspective, it worked. Many of the ContractZoo buyers were not very sophisticated and bought almost exclusively on look and feel. James was in the right place at the right time, adding customers in the aftermath of the ChatGPT release in late 2022.

But in the quiet hours after the engineering team went home, James would sometimes open a terminal, paste a contract directly into Claude, or GPT, or Gemini—whichever model he happened to be testing that week—with a well-crafted prompt, and compare the raw output against what ContractZoo produced. The differences were cosmetic. Formatting. Structure. Presentation. The analytical substance was all too often identical or nearly identical. His “proprietary legal intelligence™ layer” was, if he was being ruthlessly honest with himself, a polished prompt library with a polished front end. He had built a beautiful house on land he did not own. And the landlords, Anthropic, OpenAI, Google, and whoever came next, were getting better every quarter.

Every model release closed the gap between what ContractZoo offered and what a user could get by talking directly to the model. Someday—and James could not tell whether “someday” meant zero years or five—the models would just do the job themselves. A general counsel could open Claude or Gemini, upload a contract, ask for an analysis, and get output indistinguishable from what ContractZoo produced. No subscription. No platform. No ContractZoo.

He had arrived at a fortunate moment, when the models leaped from barely functional to genuinely capable, and he had moved fast enough to capture the wave. But riding a wave was not the same as owning the ocean. The thought that kept him up at night was that he had spent too much time on the interface and not enough building something genuinely proprietary: a data layer, a learning loop, a moat that did not depend on the models staying slightly less convenient than his wrapper around them.

James understood he needed an exit strategy. The problem was he had taken a lot of VC money and as such the list of potential buyers was somewhat lim-

ited. He had quite a bit of revenue but understood he was likely operating on borrowed time.



5.4 FIVE LAYERS OF THE CANDOR OS STACK

Priya spent the first two weeks after the build decision doing what she called “architecture from the ground up.” She had joined Sarah’s firm from a legal tech company where she had led engineering for three years, watching that company make every mistake she was now determined to avoid. Her previous employer had built features before architecture, bolted on integrations after the fact, and accumulated technical debt that eventually slowed development to a crawl. Priya was not going to repeat that pattern.



THE FIVE-LAYER PLATFORM ARCHITECTURE

Layer 1: Foundation Models and Inference—Commercial AI APIs (Claude, GPT). Buy, do not build. Design for model-agnosticism. Allocate inference-time compute strategically: routine tasks get fast, cheap passes; high-stakes analysis gets extended reasoning where the model thinks longer and deeper before responding.

Layer 2: Retrieval and Context—RAG or other retrieval and context aggregation system that grounds AI outputs in the firm’s accumulated knowledge. The compounding advantage.

Layer 3: Workflow Orchestration (with a Path to Agents)—Initially, structured pipelines that route work through AI and human stages as directed graphs with configurable thresholds and handoffs. The architecture is designed to evolve toward AI agents that can plan multi-step tasks, decompose complex engagements into subtasks, and execute them autonomously or with human checkpoints. Agents will access external systems through tools and standardized protocols like MCP (Model Context Protocol), connecting the models to document repositories, research databases, regulatory feeds, and internal knowledge bases. The pipeline comes first; the agents come when the models are ready.

Layer 4: Application Layer—Interfaces for lawyers (review dashboard), clients (submission portal), and administrators (monitoring console).

Layer 5: Analytics and Learning—Performance measurement, error tracking, correction capture, and feedback loops that make every engagement improve the next.

She presented her architecture to the full team—all eight of them, crowded into the firm’s small conference room in downtown Phoenix—on a Thursday afternoon. The whiteboard showed five layers, stacked like a wedding cake, each one supporting the layers above it.

“Layer one,” Priya said, pointing to the bottom. “Foundation models and inference. We are using Claude as our primary model for analysis and generation,

with GPT as a secondary for specific tasks where it performs better. We access them through APIs. We pay per token. This layer is bought, not built. It will change as the models improve. What matters is that our architecture is model-agnostic. If a better model appears next month, we can swap it in without rewriting everything above it.”

She added a note to the side of the diagram. “One thing that matters here is how much compute we allocate at inference time. Not every task deserves the same amount of thinking. A routine NDA review gets a fast, cheap pass. But a complex cross-border acquisition agreement with regulatory implications—that gets extended reasoning, where we let the model think longer and deeper before it responds. The newer models support this natively. We can tell the model to plan its analysis before executing, to consider multiple interpretations, to check its own work. The cost per query goes up, but the quality difference on hard problems is substantial. We need to be smart about when to spend that compute and when not to.”

“Layer two is retrieval and context. This is our RAG system, retrieval-augmented generation.” She saw a few blank looks and paused. “When we ask Claude to analyze a contract, we do not just hand it the document and say ‘go.’ We first search our knowledge base for relevant context: similar contracts we have reviewed before, client-specific preferences, common risk patterns in this industry, regulatory requirements that apply. We retrieve that context and pass it to the model along with the document. The model’s analysis is grounded in our accumulated knowledge, not just its general training.”

“That is why our analysis is better than ContractZoo’s,” Sarah added, the connection clicking. “They do not have our knowledge base. They have a generic one which must stay broad enough to handle the needs of a wide range of customers.”

“We are heading toward a world of mass customization. Namely, client and situation specific automated customization,” Priya said, “but ContractZoo has real limitations which will not make it easy for them to help deliver that future.”

“Exactly. And the knowledge base grows with every engagement. That is the compounding advantage.” Priya moved up the diagram. “Layer three

is workflow orchestration. This is where we define how work actually flows through the system.”

“A contract review engagement does not hit the AI as a single monolithic task. It moves through a pipeline: document ingestion, document decomposition, clause extraction, risk classification, cross-reference analysis, client-specific preference matching, draft report generation. Each step is a discrete node in the workflow. Some nodes are pure AI. Some require human review. Some are conditional—if the AI’s confidence score on a particular clause falls below a threshold, it routes to a senior reviewer instead of proceeding automatically.”

David’s eyes lit up. “That is a process map. You are building a process map into the technology.”

“Yes. And the process map is configurable. When we onboard a new client, we can adjust the workflow—add steps, change thresholds, route certain work types to specific individuals who have domain expertise. The workflow is our operating system, and it is designed to evolve.”

She paused and added a note to the side of the diagram. “I want to flag where this is headed, even though we are not building it yet. Right now, the workflows are structured pipelines: defined steps, defined routing, defined outputs. The AI executes within each step, but the plan is ours. Someday—and I think within a year—the models will be good enough that we can move to an agent-based architecture. Agents that can receive a complex task, decompose it into sub-tasks, figure out what information they need, go get it through standardized tool connections, something called MCP, Model Context Protocol, and execute the analysis with the ability to revise their approach when they hit something unexpected.”

She drew a box on the whiteboard and labeled it *Agent (Future)*. “That is where the real power. An agent does not just execute a plan someone else wrote. It builds its own plan on the fly. ‘It plans the work and works the plan.’ But the models are not reliable enough for that yet in production, especially when malpractice liability is on the line. So we start with structured pipelines and we design the architecture so that when the models are ready, we can drop agents into the same orchestration layer without rebuilding everything else.”

“So it thinks before it acts,” David said. “Eventually.”

“Eventually. For now, we do the thinking. The system does the executing. That is version one.” Priya wrote *MCP* on the board. “But we are going to wire up the tool connections now, MCP servers for our document store, our knowledge base, our client preference database, our regulatory feeds, so that when the agent layer is ready, the plumbing is already in place. We are building the roads before we have self-driving cars.”

“Layer four is the application layer, the interfaces people actually use. Our lawyers see a review dashboard where AI-generated analyses queue for verification. Clients see a portal where they can submit documents, track progress, and download deliverables. I see an admin console where I can monitor system performance, adjust configurations, and deploy updates.” Priya turned back to the group. “Layer five is analytics and learning. Every engagement generates data. David, this is your layer as much as mine.”

David opened his notebook. He had been preparing for this moment. “Every time a reviewer accepts, modifies, or rejects an AI-generated analysis, that decision is captured. Not just the outcome, accept or reject, but the nature of the correction. Did the AI miss a clause? Did it misclassify a risk? Did it fail to account for a client-specific preference? Did it flag something as high-risk that was actually standard?” He looked around the room. “That data is structured. It is tagged by clause type, contract type, client, reviewer, and date. Over time, it tells us exactly where a given model is strong and where it is weak. It tells us which reviewers catch which kinds of errors. It tells us which client contexts are hardest for the AI to handle. And it feeds back into layer two, the knowledge base, to make the next analysis better.”

The room was quiet for a moment. Sarah’s eyes settled on the architecture diagram, and she felt something shift. This was not a collection of tools anymore. It was slowly becoming a platform. And the platform was designed to learn.

“How long to build this?” she asked.

Priya did not hesitate. “MVP in three to six months. Production-grade in six to twelve. But I need resources.”

One of the lawyers, Joshua, who had been with the firm since the beginning, raised his hand. “Can I ask a dumb question?”

“No dumb questions,” Sarah said, though she suspected from David’s expression that he was mentally cataloging this as a process improvement opportunity.

“Why do we need all five layers right now? We are eight people. Can we not just...keep using Claude and fix the Google Sheet problem?”

Priya nodded. “Fair question. You can absolutely keep doing what you are doing today. You actually have to keep doing it in the short term while we do the build out. But if we do not move beyond the approach used today, here is what happens as you scale. You grow to twelve people. Now your Google Sheet has six editors, and the data conflicts multiply. You take on a client who needs two hundred contracts reviewed in three weeks. Your current workflow cannot handle the volume because every document requires manual routing. A reviewer catches an error that the AI has been making consistently for months, but there is no mechanism to prevent it from happening again. Each of these problems is solvable individually. But solving them individually, ad hoc, without architecture, is how you end up with the duct-tape system you have now—except bigger and more fragile.”

Joshua nodded slowly. “So we are building the road before the traffic arrives.”

“Exactly. It is much cheaper to build the road before the traffic than to repave while cars are driving on it.”

The meeting ended with Priya assigning the first sprint’s tasks. Sarah lingered, looking at the whiteboard. Five layers. Three months to MVP. A budget request she had not yet seen the numbers for. She could feel the weight of the decision settling onto her shoulders, the familiar tension between what the firm needed to become and what it could afford to build right now.

5.5 YOUR BUDGET REVEALS WHAT YOU ARE

The conversation Sarah had been dreading arrived two days later, on a Saturday morning board call. Alex Hawthorne dialed in from Los Altos. Sarah sat in her home office with Priya’s budget proposal open on her screen. The number at the bottom made her stomach tighten.

Priya was asking for twenty-five percent of revenue allocated to technology: engineering salaries, cloud infrastructure, API costs, tooling. For a firm generating roughly \$1.3 million in annual run-rate revenue, that meant roughly \$325,000 per year dedicated to technology. But that was not going to be enough to build this out. The engineering hires Priya requested would consume a significant share of their remaining seed capital (unless they could increase revenue in the meantime).

“Walk me through your objection,” Alex said, his voice calm through the speaker.

“My objection is that twenty-five percent of revenue on technology is insane for a law firm,” Sarah said. “Traditional firms spend two to five percent. Even the most aggressive ones I have seen spend maybe eight. Priya is asking for five times that. And we already pay her quite a bit of money and she has a very aggressive equity vesting schedule.”

“Okay sure that is what it takes to get genuine technical talent,” said Alex. “But what does Priya say?”

“She says it is the minimum to build the platform properly. She says anything less and we are back to duct tape. She says the firms that underfund technology end up with systems that cannot scale, and then they spend twice as much fixing it later.”

There was a pause on the line. When Alex spoke, his tone had shifted—the casual warmth replaced by something more direct, the voice of someone who had spent decades evaluating businesses and was about to say something he considered important.

“Sarah, I am going to tell you something, and I need you to actually hear it, not just acknowledge it.” He paused. “You are not a law firm that uses technology. You are a technology company that practices law. Fund it accordingly.” After a few seconds of silence he added “if anything you need fewer lawyers.”

Sarah opened her mouth to respond, then closed it. The sentence landed with a weight she had not expected. She had been thinking of the technology investment as a cost center, a necessary expense that supported the firm’s real business of delivering legal services. Alex was telling her to invert the framing entirely. The technology was the business. Legal services was the application.

“The traditional firm benchmarks are irrelevant,” Alex continued. “A traditional firm spends a single digit percentage on technology because technology is not their competitive advantage. Their competitive advantage is partner relationships and associate labor. Your competitive advantage is your AI platform and your data. Seen in this light, twenty-five actually seems kinda low. If you underfund the thing that differentiates you, you are just another small law firm with a crappy chatbot.”

After the meeting, Sarah recounted a (sanitized) version of the Alex discussion with Priya. “Listen, I have seen this play out at my last company. They tried to build a technology platform on a somewhat skimpy technology budget. The result was a system that worked well enough for demos but could not handle production load. When clients started sending real volume, the system buckled. They lost three customers in one quarter because deliverables were late and quality was inconsistent. By the time they decided to invest properly, they had already damaged their reputation. We have to have great tech here that is core to what Candor offers. In my opinion, there was no reason to start this firm and raise money if you were not going to put a major bet on tech.”

Sarah pulled up the budget spreadsheet. Twenty-five percent of revenue today. She ran the numbers forward. If revenue grew as projected, the absolute dollar amount would increase, but as a percentage it would gradually decline as the platform matured and required less new development. By year three, Priya projected technology would consume fifteen to eighteen percent of revenue—still high by traditional standards, but within the range that the analytical models suggested for a growth-stage AI-native firm.

“What if revenue does not grow as projected?” Sarah asked.

“Then we have a bigger problem than the technology budget,” Priya said. “If revenue stalls, the technology spend did not cause it. Underfunding the platform is more likely to cause it.”

Sarah studied the numbers one more time. Then she closed the spreadsheet.

“Okay. Twenty-five percent and we are going to spend a bunch of our capital to do the build. So this is a serious commitment and reflects our view that we are indeed a tech company. But I want monthly reporting on what that money is producing. Not just activity. Output. Features shipped, system performance,

quality metrics. David, I want you tracking the return on this investment the same way you track everything else.”

David, who had been quietly following the conversation, nodded. “I will build it into the operational dashboard. Technology spend will have the same accountability as every other line item.”

One other thing that Alex had mentioned stuck with Sarah. Priya needed to use AI to build AI: more vibe coding and less hand-stitching by needlepoint. Sarah could see that Priya was somewhat reluctant to the idea. Alex had been insistent. He said that his best companies were getting huge gains by using tools such as Claude Code and similar AI-assisted development environments.

Old habits die hard, and Priya was processing the same sort of refactoring that every lawyer who worked at Candor also had to face. The world was different now and she had to adapt. So Priya said the right thing while she slowly tried to do the right thing. “Ok yes I have been avoiding all of this but we are a tech company and we have to...what is the phrase ‘eat our own dog food.’”

“Good,” Sarah said. “Alex’s view is that firms that win in this space are not the ones that spend the most on technology. They are the ones that spend it on the right things. Priya’s architecture, the five layers, that is the right thing so long as it is built under the best possible economics.”

5.5.1 Technology Investment for AI-Native Firms

Sarah’s budget discussion illustrates a tension that every AI-native firm must resolve. Traditional professional services firms spend a single digit percent of revenue on technology because technology supports their operations without driving differentiation. AI-native firms operate on different economics entirely.

Early-stage AI-native firms are likely to allocate twenty-five to fifty percent of revenue to technology: engineering talent, cloud infrastructure, API costs, and development tools. This level reflects the reality that the technology platform is the firm’s primary competitive asset, not a support function. As the platform matures and revenue scales, the percentage gradually declines to a more modest

percent of revenue, though the absolute investment continues to grow. An AI-native firm simply spends a much higher percentage of revenue on technology than the traditional service provider.

It is easy for folks to call themselves AI-native. An increasing number of incumbents, ALSPs and other service providers, claim to have been AI-transformed. The level of technology investment, not the marketing, is a good bellwether for the legitimacy of such claims.

The allocation within the technology budget matters as much as the total. At scale, engineering talent typically consumes thirty-five to forty-five percent of the technology budget, the people who build and maintain the platform. AI models and compute costs account for twenty-five to thirty-five percent, the foundation model API calls and infrastructure that power the analysis. The remainder covers third-party vendor tools, security, and infrastructure.

The key insight from Sarah's board call, Alex's reframing of the firm as "a technology company that practices law," is not merely rhetorical. It determines how the firm allocates capital, how it prioritizes investments, and how it measures return. A law firm that uses technology evaluates tech spending as overhead to be minimized. A technology company that practices law evaluates tech spending as R&D to be refined for competitive advantage.



5.6 DATA AS THE FIRM'S MOST VALUABLE ASSET

David Park had been thinking about data since before Priya arrived. His background in Lean and Six Sigma had trained him to see every process as a source of measurement, and every measurement as an opportunity for improvement. In previous roles, he had seen the deep process-oriented knowledge that lived in people's heads never reach a system where it could be memorialized. When the people left, the knowledge left with them.

He was determined not to let that happen at Sarah's firm. But the challenge was harder than he had expected. AI systems generated a different kind

of institutional knowledge—not the war stories and relationship insights that departing partners took with them, but structured data about what the AI got right, what it got wrong, and what human reviewers did to fix it. That data could become the firm’s most durable competitive asset if it were properly captured, organized and operationalized into systematized improvement protocols executed by a combination of humans and agents.

David presented his data strategy to Sarah and Priya over lunch one Tuesday, spreading printouts across the table at Phx-Pho, the team’s favorite Vietnamese restaurant located just a few blocks from the office.

“Every engagement we run generates three categories of data,” he said. “Input data: the client documents, the engagement parameters, the specific questions the client wants answered. Output data: the AI-generated analysis, the human-reviewed deliverable, the final product we send to the client. And correction data: the delta between what the AI produced and what the human reviewer changed before it went to the client.”

“The correction data is the gold,” Priya said immediately.

“Exactly. The correction data tells us where the AI fails. Not in the abstract—specifically. On this clause type, in this contract context, for this kind of client, the AI tends to miss this particular risk. Or it flags something as high-risk when industry practice treats it as standard.” David held up a printout showing a spreadsheet template. “I want every reviewer to log their corrections in a structured format. Not free-text notes in the margin of a document. Structured fields: clause type, error type, severity, correction made, reasoning.”

Sarah frowned. “That adds time to every review. Our lawyers are already stretched.”

“Two minutes per correction, on average. I timed it.” David met her eyes. “Two minutes per correction in exchange for a data asset that makes the AI better on every subsequent engagement. The ROI on that two minutes compounds indefinitely.”

He leaned forward. “The key is to make distinctions between customer-specific requirements and generalizable insights. Both are useful but the key is to apply the right lens to each problem. The right blend of broad general insights

combined with customer-specific customization only comes from a two-tier flywheel of general and specific data feedback loops.”

It was a grand vision and Priya was already thinking about the technical implementation. “We can build the correction capture into the review interface. When a reviewer modifies an AI output, the system automatically records the before and after. The reviewer just needs to tag the error type and add a brief note on the reasoning. We can make the tagging a dropdown menu: fast, structured, consistent.”

“And the corrections feed back into our overall analysis pipeline?” Sarah asked.

“Yes. When the AI encounters a similar clause in a future engagement, it retrieves not just the general knowledge but the specific corrections from previous reviews for this client and this client’s preferences. Over time, the error rate on common clause types should decline measurably.” Priya paused. “But RAG, retrieval augmented generation, is a hack really; we are going to use this to set up our future agentic-driven offerings.”

Sarah said, “The real benefit of being a service provider today is that we can overlay any shortcomings in our platform with human expertise. We deliver high-fidelity output while creating the learning loop that allows us to improve every day.”

David said. “This is the learning loop that ContractZoo does not have. Their corrections stay in their customers’ heads. Ours stay in our system.”

David had one more point. “I want to track correction rates by reviewer, too. Not to punish anyone—to understand variance. If one reviewer catches errors that others miss, I want to know why. Is it domain expertise? Is it a different review methodology? Whatever it is, we can potentially codify it and improve the whole team’s performance.”

Sarah watched the two of them—the operations mind and the engineering mind, converging on the same insight from different directions. Data was not a byproduct of their work. It was the point of their work. Every contract reviewed, every correction logged, every client preference captured, all of it fed a system that grew smarter with use. A competitor could license the same models, hire the same caliber of lawyers, even copy their workflow design. But they could not so

easily replicate years of accumulated correction data, client-specific knowledge, and performance analytics.

DATA AS COMPETITIVE MOAT

The models are commodities. The workflows can be copied. But proprietary data accumulated over time creates an advantage that compounds. Each engagement adds correction data, client preferences, and performance analytics. A firm with years of structured correction data has a capability edge that no new entrant can replicate quickly. The moat grows through volume, quality, uniqueness, and network effects—and it only exists if you design for data capture from day one.

“This is the data moat Maya talked about,” Sarah said, half to herself.

“What moat?” David asked.

“A friend of mine, an investment banker, told me early on that the most defensible competitive advantage an AI-native firm can build is its data. The models are commodities. The workflows can be copied. But proprietary data accumulated over time creates an advantage that compounds. She was right. I just did not fully understand how right she was until now.”

5.6.1 *The Data Strategy in Practice*

Candor implemented David’s structured data capture within three weeks of the conversation. The system was simple at first: correction fields embedded in the review interface, with dropdown menus for error classification and free-text fields for reasoning. Priya’s engineering team built a pipeline that ingested the correction data nightly, updating the RAG knowledge base and generating weekly analytics reports.

The results were measurable within the first month. The AI’s accuracy on standard indemnification clauses, one of the most common clause types in their contract review work, improved from eighty-two percent to eighty-nine

percent after incorporating corrections from the first set of engagements. Error rates on limitation-of-liability provisions dropped by a third. The improvement was not uniform. Some clause types proved stubbornly resistant to correction, particularly those involving multi-party interactions where context from distant parts of the contract mattered. But the trend was unmistakable.

David tracked the metrics with the rigor of someone who had spent a career measuring process improvement. He built a dashboard that showed correction rates by clause type, by client, by reviewer, and by time period. The dashboard became a centerpiece of the firm's weekly operations meeting, a visual representation of the firm getting smarter with every engagement.

The data strategy also raised governance questions that Sarah had not anticipated. When Priya and David proposed feeding client document data into the RAG knowledge base, Sarah's legal instincts flared.

"We need client consent," she said flatly during the governance discussion. "We called this firm Candor for a reason. Every engagement letter from here forward needs to include a clear, specific authorization for us to use de-identified engagement data to improve our AI systems. Existing clients need to be contacted and asked to opt in. No assumptions, no implied consent."

Priya pushed back gently. "If we exclude existing client data, the knowledge base starts nearly empty. The learning loop is slower."

"I understand the trade-off. But we are a law firm. Client confidentiality is not something we work around—it is a constraint we design within." Sarah thought about some of the data failures from across various industries. Several had involved data handling practices that, while technically legal, eroded client trust. "If a client finds out we used their contract data to train a system without telling them, we do not just lose that client. We lose our reputation. And for a firm called Candor, that would be more than embarrassing—it would be fatal."

David supported Sarah's position. "We can design the consent process to be lightweight. A clause in the engagement letter, a clear explanation of what we do and do not do with the data, an easy opt-out. Most clients will say yes if you explain the benefit: their future work gets better because the system learned from their past work. To those who do not consent, we can eventually consider either differential pricing or simply move on from them as clients."

In the end, they achieved an eighty-five percent opt-in rate among existing clients. The remaining fifteen percent had their data excluded from the knowledge base entirely, a clean separation enforced by access controls that Priya built into the system architecture. The consent process added a small amount of friction to client onboarding, but it also became a selling point: clients appreciated that the firm was transparent about its data practices, a contrast to the opacity they encountered with most legal tech vendors. The fact that lawyers and not some random tech person were the ones in charge of data governance was a major advantage that the AI-native law firm had over pure product companies.



5.7 THE PLATFORM TAKES SHAPE

By the start of Priya’s fourth month, January 2026, the platform had a name. Internally they called it Candor OS. The MVP was functional, though production-grade maturity was still three months away, on track for spring.

The name was deliberate. David had been calling their workflow “the operating system” since the early weeks, and Priya made it literal. “An operating system sits between the hardware and the applications,” she told the team. “It manages resources, handles complexity, and lets everything else run. That is what this platform does. It sits between the AI models and the legal work, orchestrating everything so the lawyers can focus on judgment, not plumbing.”



She paused. “And if we build it right, the OS becomes the foundation for something bigger. Right now, humans initiate every workflow. Eventually, some workflows will run fully autonomously, agents handling routine matters end to end, with humans reviewing outputs rather than driving the process. For higher-complexity work, the agents assist rather than replace, surfacing research,

flagging risks, drafting analyses that lawyers refine and approve. Both modes run on the same OS. That is the whole point. You do not get truly agentic without the operating system first.”

The platform had a shape that matched her five-layer architecture. The foundation model layer was stable, with Claude handling primary analysis and generation tasks through well-structured API calls with retry logic and fallback routing. The retrieval layer had been rebuilt from scratch, replacing the ad hoc collection of Notion pages and Google Docs with a proper vector database using domain-aware chunking that split contracts at clause boundaries rather than arbitrary token limits.

The workflow orchestration layer was where Priya’s engineering instincts showed most clearly. She had designed each workflow as a directed graph: a series of nodes connected by edges, where each node represented a discrete processing step and each edge represented a routing decision. A contract review workflow might have fifteen nodes: document ingestion, format normalization, clause extraction, clause classification, risk scoring, cross-reference detection, client preference matching, draft annotation, confidence scoring, routing (automated path for high confidence, human review for low confidence), human verification, correction capture, report generation, quality check, and client delivery.

“The beauty of the graph architecture,” Priya told Sarah during a late-evening working session, “is that we can modify any node independently. If we find a better approach to clause extraction, we swap out that node without touching the rest of the pipeline. If a client needs an additional analysis step—say, regulatory compliance screening for their industry—we add a node. The workflow grows with the business.”

Sarah watched Priya work and thought about the gap between this and what she had been doing six months ago. Back then, “contract review” meant opening a document in Claude, writing a prompt, reading the output, editing it into a memo, and emailing it to the client. Even with “AI” in the mix, each engagement was a craft production operation. Skilled, personal, and utterly non-scalable. Bottlenecked by the whims and limits of human reasoners prompting a system like a snake charmer playing a flute.

Now the same work flowed through a system that learned, adapted, and improved with each repetition. The humans were still essential, Priya was emphatic about that, but they were essential for different reasons. Not for the reading and summarizing that consumed most of a traditional associate's time, but for the judgment that no system could replicate: understanding the client's strategic context, weighing risks that required business judgment, catching the subtle errors that pattern-matching alone would miss.

"We are not replacing lawyers," Priya said, as if reading Sarah's thoughts. "We are replacing the parts of lawyering that arguably never required a law degree in the first place."

5.7.1 A Litany of Technical Design Decisions

Priya made several design decisions that reflected hard-won lessons from her previous role. The chunking strategy split documents at semantic boundaries (section headings, clause breaks, paragraph transitions) rather than at fixed token counts. This preserved the logical structure of legal documents, where a single clause might span multiple paragraphs and breaking it mid-thought would destroy the context a model needed to analyze it correctly.

The search approach combined semantic retrieval with keyword matching. Legal documents contained defined terms, statutory references, and precise language where exact-match search outperformed the fuzzier results of embedding-based retrieval. A hybrid approach ensured that when a contract defined "Change of Control" with specific meaning, the system found that definition reliably, not just passages that were semantically similar to the concept of change of control.

Priya also implemented a reranking step after initial retrieval. The first-pass search returned candidate chunks ranked by relevance. A second model then rescored those candidates with more precise relevance judgments, filtering out passages that were topically related but not actually useful for the specific analysis task. The reranking added a fraction of a second of latency but meaningfully improved the quality of context passed to the generation model.

“Garbage in, garbage out applies to RAG as much as anything else,” Priya told the team. “The generation model is only as good as the context we give it. If we retrieve irrelevant passages, the model will confidently incorporate irrelevant information into its analysis. The retrieval layer is not a search engine—it is a knowledge curator.”

5.7.2 *A Live Fire Drill*

The platform’s first real stress test came six weeks after Priya’s architecture presentation. A mid-market logistics company, one of the firm’s newer clients, needed 427 vendor contracts reviewed before a corporate restructuring, with a three-week deadline. Under the old workflow, Sarah would have declined the engagement or negotiated a longer timeline. Even though they were still in Beta, Priya believed they could handle it with Candor OS.

The engagement exposed both the platform’s strengths and its gaps. The document ingestion pipeline processed all 427 contracts in under two hours—a task that would have taken a paralegal a full week of scanning, organizing, and filing. The clause extraction and risk classification layers performed well on standard provisions, producing first-draft analyses that reviewers could verify in fifteen to twenty minutes per contract rather than the hour-plus a cold read would require.

But the cross-reference detection failed on a particular pattern common in logistics contracts: force majeure provisions that interacted with separate service-level agreements referenced by exhibit number. The AI identified the force majeure clauses and flagged the SLA references independently, but it did not connect the two. It could not understand that the practical effect of the force majeure provision depended entirely on the SLA terms in a separate exhibit that the system had indexed as a different document.

Sarah’s senior reviewer caught the gap on the fourth contract and immediately flagged it. David’s correction pipeline captured the error type, and Priya’s team adjusted the retrieval logic within two days to link cross-referenced exhibits during the ingestion phase. By the fortieth contract, the system was handling the force majeure-SLA interaction correctly. By the hundredth, it was catching variations that the reviewers had not yet encountered.

“That is the learning loop in action,” David said at the post-engagement retrospective. “One error, caught by a human, corrected in the system, never repeated. At a traditional firm, that error would have been caught by one associate, mentioned over coffee to another associate, and forgotten within a week. The institutional memory was the partner or associate. Our institutional memory is the platform.”

The logistics engagement was delivered on time, under budget, and with a client satisfaction score that Sarah found gratifying. More important than the client outcome was what the engagement taught the team about the platform’s capability trajectory. Each large engagement was not just revenue—it was a stress test that revealed weaknesses, generated correction data, and left the system measurably better than it had been before. This kaizen loop was slowly but surely becoming an enterprise value creation machine.

Sarah tracked the firm’s progress against the AI-Native Maturity Model she had developed in her early planning. A year ago, her firm had been at Level 2, using AI tools but without integrated workflows or systematic data capture. The infrastructure Priya was building would push them to Level 3 and, within the year, toward Level 4.

The progression was tangible. At Level 2, each engagement had been a standalone effort, dependent on the individual lawyer’s skill at prompting an LLM and the team’s collective memory of past work. At Level 3, engagements flowed through defined workflows that incorporated the firm’s accumulated knowledge. The quality was more consistent, the turnaround faster, and the capacity meaningfully higher.

Level 4, AI-native operations, was the target. At that level, the AI would handle the majority of analytical work, with humans supervising, verifying, and adding judgment. The firm’s competitive advantage would come from its platform and data, not from the individual skills of its lawyers. That did not mean the lawyers were less important. It meant they were important for different, higher-value reasons.

But Sarah was beginning to realize the framing needed refinement. The platform’s most valuable component was not the data or the retrieval logic—it was the encoded judgment. The analytical frameworks Priya had embedded in

the prompt architecture, the risk-weighting criteria that lawyers had refined over dozens of engagements, these were the firm's real intellectual property.

Every competent firm in a given practice area had roughly the same templates and forms: the NDA, the stock purchase agreement, the vendor services agreement. These were commodity inputs. What differentiated Candor's output was the instruction set that told the AI how to think about the work: what to flag, what to ignore, how to weigh competing considerations, what level of confidence to require before making a recommendation. Those playbook-style instructions encoded individual professional judgment into a form that could be transferred across the entire organization. A new hire did not need years of apprenticeship to internalize Sarah's analytical framework for regulatory compliance review; the framework was in the system, applied automatically, refined with every engagement. The template was a commodity. The judgment encoded in the instructions was the moat.

David framed it in operational terms during a Friday retrospective. "At Level 2, we were artisans. Each piece of work was handcrafted. Beautiful, but slow and variable. At Level 3, we are a workshop—defined processes, shared tools, consistent quality. Level 4 is a factory with artisan finishing. The system handles production. The humans handle the work that requires human judgment. The combination is faster, more consistent, and more scalable than either alone."

"I am not sure our lawyers love the factory metaphor," Sarah said.

"They do not have to love it. They have to understand that the alternative is being replaced by someone else's factory." David paused. "The firms that resist systematization are not preserving craftsmanship. They are preserving inefficiency. And their clients will eventually understand this."

Sarah noted the distinction. The maturity progression was not just about technology capability—it was about organizational capability. Moving from Level 2 to Level 3 required the technology, yes, but it also required the team to change how they thought about their work.

The lawyers had to accept that verification was as valuable as drafting. The operations team had to embrace structured data capture as part of their daily routine, not an administrative burden. The engineers had to build for the firm's actual needs, not for technical elegance. Each level demanded a cultural shift

alongside the technical one, and the cultural shift was often the harder of the two.

Firms that attempted to skip levels—installing Level 4 technology on a Level 2 organization—typically failed. The technology outran the team’s ability to use it effectively. Conversely, firms with Level 4 organizational maturity but Level 2 technology were constrained by their tools but positioned to leap forward once the technology caught up. Sarah’s advantage was that she was building both simultaneously, with Priya driving the technology maturity and David driving the organizational maturity in parallel.



5.8 BASE CAMP BUT NOT THE SUMMIT

Sarah hit the Echo Canyon trailhead at six on a Saturday morning, before the heat turned Camelback Mountain into a furnace. She had not been a hiker before her trip to Palo Alto. But watching the Sand Hill Road crowd treat their morning trails like moving conference rooms had made her realize she had spent a decade feeding her mind while ignoring her body. Now it was a weekend ritual—usually solo, always early.

The trail was steep from the start. Her legs burned within the first quarter mile, the red rock scrambles demanding hands and knees in places. Other hikers passed her going up, some faster, some slower. A pair of women in matching Lululemon moved with the easy cadence of people who did this every day. A man in his fifties labored ahead of her, stopping every few minutes to catch his breath but never turning back. She related to him more than team Lululemon.

It had been about fourteen months since she started Candor. She had eight people, a growing client list, a platform architecture that slowly was starting to feel real, and a data strategy that would compound over time.

The ContractZoo pilot felt like a lifetime ago. She had been tempted—genuinely tempted—to buy a solution and skip the hard work of building one. The UI was polished. The ContractZoo demo was far more impressive than

Candor OS. But as Priya reminded her, “the UI just needs to be functional for an internally used tool. We are not selling to random lawyers where the UI needs to be flawless.”

The prospect of spending months and hundreds of thousands of dollars building custom technology had once been daunting for someone whose entire career had been spent practicing law. Priya, it turned out, was up to the task.

But the pilot had revealed something important: the difference between a product and a platform. ContractZoo was a product, a packaged solution designed for the average legal team. What Sarah was building was a platform: a system designed to learn, adapt, and improve based on the specific knowledge and experience of her firm. The product would always be limited by its vendor’s roadmap and its generic templates. The platform would grow with every engagement, shaped by the firm’s own data and the corrections of its own reviewers.

She paused at a ledge about two-thirds up, breathing hard, and looked out over the Valley. Phoenix sprawled below in every direction: the grid of streets, the clusters of development, the brown desert pushing in at the edges. This was the city where she had chosen to build. Not San Francisco, not New York, not the places where legal innovation was supposed to happen. Phoenix.

Lower cost of living, lower overhead, and a local client base of mid-market companies that needed sophisticated legal work but were generally more cost sensitive and practical. The geography was a strategic choice, even if it sometimes felt like exile.

Alex’s words echoed: *You are not a law firm that uses technology. You are a technology company that practices law.*

She had resisted that framing at first. It felt like an identity she had not earned—a lawyer calling herself a technologist. But watching Priya build the platform and David build the data infrastructure, she understood what Alex meant.

The technology was not a tool the firm used. It was the firm’s core asset. The legal expertise mattered enormously—without it, the AI’s output would be unverifiable, the corrections meaningless, the client relationships impossible. But the legal expertise was expressed *through* the platform, not alongside

it. The platform was how the firm's knowledge became durable, scalable, and compounding.

Her phone buzzed in her pocket. She pulled it out. A text from Priya:

“RAG layer v2 just passed integration tests. Retrieval accuracy up 23% on the benchmark set. Also, David's correction pipeline is feeding data cleanly into the knowledge base. First full learning loop is operational.”

Sarah smiled, sweat dripping onto the screen. The first full learning loop. An engagement would generate corrections. The corrections would update the knowledge base. The updated knowledge base would improve the next engagement. The cycle would repeat, and the firm would get smarter each time.

It was not dramatic. No champagne corks, no celebration. Just a text message on a Saturday morning on the side of a mountain, confirming that a system designed to learn had started learning.

She typed back, thumbs clumsy from the climb:

“Good. Let's see what the correction rate looks like after fifty more engagements. And tell David his data capture spreadsheet is the most valuable thing in this firm.”

Priya's reply came immediately:

“Don't tell him that. He's already insufferable about process documentation.”

Sarah laughed and put the phone away. She looked up at the summit—still a hard scramble above her. The view from here was already impressive. You could see how far you had come. But you could also see how much mountain was left.

That was the firm, too. The client pipeline was growing. The platform was taking shape. The data was accumulating. Somewhere in the intersection of legal expertise, AI capability, and operational discipline, a new kind of firm was emerging. But they were at base camp, not the summit. There was more to build. There would always be more to build.

Monday would bring a client onboarding, a sprint review with Priya's engineers, David's weekly operations meeting where he would present the latest correction rate trends and push the lawyers, diplomatically but relentlessly, to be more consistent in their tagging. There would be vendor calls to evaluate a new

embedding model, a security audit to prepare for, and the ever-present question of when to hire the next engineer versus the next lawyer.

She started climbing again. The man in his fifties was still ahead of her, still stopping, still not turning back. The summit was up there somewhere above the rocks and the scrub. She did not know what the firm would look like in three years, but she knew it would be better than it was today, because every engagement between now and then would leave it smarter than the one before.

That was the technology thesis, reduced to its simplest form: build a system that learns from its own work, and then do a lot of work.

COMING NEXT WEEK

Chapter 6
Land, Expand, Compound

The platform is humming. The data flywheel is turning. But Sarah still does not have an answer to the most fundamental question: what kind of firm is she actually building? Replace—let AI do the work while lawyers supervise? Optimize—keep lawyers at the center, just make them faster? Chapter 6 dismantles that false binary, introduces the AI-Era Positioning Matrix, and forces a strategic choice that will define every hire, every investment, and every client relationship going forward.

Subscribe at <https://theainativefirm.com> to receive new chapters as they release.

ABOUT THE AUTHORS

Daniel Martin Katz, PhD, JD is Professor at Illinois Tech–Chicago Kent College of Law and Academic Director of the Bucerus Center for Legal Technology & Data Science. Named by the *Financial Times* as one of the top 20 legal market shapers of the past twenty years, his research focuses on legal analytics, legal technology, and the future of the legal profession.

Michael J. Bommarito II, MSE, MA is a serial entrepreneur, researcher, and adjunct professor with twenty-five years of industry experience. He has founded and led multiple companies at the intersection of artificial intelligence and professional services.

Jillian Bommarito, CPA, CIPP/US/E is an advisor, risk and governance expert, and one of the first certified AI auditors in the world. She brings deep expertise in compliance, privacy, and the governance challenges of deploying AI in regulated industries.

Get the Full Book

<https://theainativefirm.com>