

# The Data Bottleneck: Why Andreessen Horowitz Bet \$30M on Protege

JAN 11, 2026



Share

## Table of Contents

The Exhaustion Problem

Why Travis May Built This Again

The a16z Investment Thesis

What Protege Actually Does

Why This Team Can Execute

Economic Realignment

What This Means for Builders

## Abstract

The AI industry faces a fundamental constraint that compute and architecture improvements cannot solve: access to high-quality, real-world data. Public datasets have been exhausted, synthetic data has hit scaling limits, and the vast majority of valuable data remains locked in private systems across healthcare, enterprises, and operational environments. Andreessen Horowitz's decision to lead a \$30M Series extension in Protege represents a thesis that data infrastructure will be as foundational to AI as cloud infrastructure was to SaaS. The company, founded in 2017 by Travis May and Bobby Samuels, is building the platform that connects AI builders with massive, multimodal datasets across healthcare, video, audio, motion capture,

and other domains. May previously founded and led both LiveRamp and Datavar major exits, while Samuels brings operational expertise from leadership roles at companies. This essay examines why a16z believes data access is the critical bottleneck for AI advancement, why this specific team has unique advantages to it, and what the emergence of functional data infrastructure means for AI development economics.

## The Exhaustion Problem

The progression of language models from GPT-2 to GPT-4 and beyond tells a clear story about the role of data in AI advancement. Early gains came from better architectures and more compute. Transformers beat LSTMs. Scaling laws held. 100 parameters plus more GPUs equaled better performance. But somewhere around 2023, the easy wins from architecture and compute started running into a hard wall. Not because the models could not scale further, but because the training data could not.

Common Crawl has been scraped to death. Reddit threads from 2012 have been ingested a dozen times over. GitHub repositories are exhausted. Wikipedia exists in every major model's training corpus. The entire public internet, which seemed infinite when these projects started, turns out to be finite and largely consumed. Synthetic data generation helps at the margins but cannot replace real-world complexity. Models trained primarily on synthetic data tend to collapse into repetitive patterns and hallucinate in predictable ways when faced with novel situations.

The problem extends beyond text. Computer vision models need diverse, high-quality labeled images and videos that capture edge cases, rare events, and unusual lighting conditions. Audio models need clean recordings across accents, environments, and acoustic conditions. Robotics and embodied AI need sensor data from physical environments. Medical AI needs patient outcomes across diverse populations and treatment contexts. All of this data exists, but almost none of it is publicly available or easily accessible.

Meanwhile, model architectures are converging. The difference between leading frontier models has less to do with fundamental architectural innovations and more with training data quality, instruction tuning datasets, and RLHF approaches. When Anthropic releases a new Claude variant or Google ships an updated Gemini model, the competitive advantage often comes down to what data they trained on or how they curated it, not whether they invented a novel attention mechanism.

This creates an uncomfortable reality for AI builders. The next 10x improvement in model capability will not come primarily from buying more H100s or hiring more researchers. It will come from getting access to better training data. Specifically, world data that captures the messy, multimodal, high-stakes environments where systems will actually operate. The internet represents maybe 5% of the world's total data. The other 95% sits in hospitals, enterprises, research labs, media archives, operational systems. Unlocking that data is the problem Protege is solving.

## **Why Travis May Built This Again**

Travis May has spent nearly two decades building data infrastructure companies. Protege represents his third major swing at solving data fragmentation problems; his track record is about as good as it gets in enterprise data, with two successful companies already under his belt before starting Protege at age 37.

May co-founded LiveRamp in 2011 with Auren Hoffman, initially joining as VP of Product before becoming CEO. The company, originally called Rapleaf, built identity resolution infrastructure for marketing and advertising, becoming the dominant platform for how brands connected customer data across different systems while maintaining privacy. LiveRamp was acquired by Acxiom for \$310M in 2014, later spun out as an independent public company in 2018, and at its peak was processing data connections for basically every major brand and publisher.

After LiveRamp, May reconnected with Vivek Ramaswamy, a friend from Harvard where they had co-founded Campus Venture Network together in 2007. Ramaswamy had founded Roivant Sciences and wanted to rebuild LiveRamp's data connectivity capabilities but for healthcare. In 2017, May co-founded Datavant with Ramaswamy.

with Roivant providing initial funding and May making a significant personal investment. May served as CEO of Datavant from 2017 to 2022, building it into a leading healthcare data infrastructure company. Datavant merged with Ciox Health in a \$7B transaction in 2021, creating what became the largest neutral health data ecosystem in the US, connecting over 2,000 hospitals, 15,000 clinics, and hundreds of other healthcare organizations. May stepped down as CEO after the merger, became President and remaining on the board.

The pattern across both LiveRamp and Datavant was identical. Build technical infrastructure to solve hard data problems that everyone faces but nobody wants to build themselves. Create a neutral platform that benefits from network effects as more parties connect to it. Focus on privacy, compliance, and trust as core product features, not afterthoughts. Charge reasonable prices so you become infrastructure rather than a vendor. Scale to the point where you are embedded in critical work and switching becomes prohibitively expensive.

What May saw with AI was that the same fundamental problem was emerging, just in a different context and at potentially much larger scale. Companies had spent a decade building massive model architectures and acquiring GPUs, but now they were hitting a wall on training data. The valuable data was fragmented, locked up, regulated, and hard to access. Every AI company was trying to solve this problem independently, negotiating one-off deals with data providers, building custom pipelines, and burning resources on undifferentiated heavy lifting. It looked exactly like the problem LiveRamp solved for marketing data and Datavant solved for healthcare data.

The difference this time was that AI represented a much larger market opportunity. Marketing data infrastructure was a multi-billion dollar market. Healthcare data infrastructure was bigger, with Datavant reaching a \$7B valuation. But AI training data infrastructure could be an order of magnitude larger because it cuts across all verticals and every geography. Every frontier AI lab needs this. Every AI application company needs this. Every enterprise building internal AI capabilities needs this. The total addressable market is basically the entire AI industry, which is projected to be worth trillions.

May also recognized that the window to build this was limited. Data infrastructure businesses have strong first-mover advantages because they benefit from network effects on both sides. The platform that signs up the most data suppliers first becomes more attractive to AI builders. The platform with the most AI builder demand becomes more attractive to data suppliers. Once that flywheel starts spinning, it becomes hard to displace. LiveRamp won because it got to critical mass first in marketing. Datavant won because it got to critical mass first in health data. Prot has the opportunity to win in AI data infrastructure, but only if they move fast.

In 2024, May co-founded Protege with Bobby Samuels, who he had worked with both LiveRamp and Datavant. Samuels had been at LiveRamp in partnerships role from 2014 to 2016, then joined Datavant in 2020, eventually becoming General Manager of Privacy Hub before leaving in 2023. The two brought in Engy Ziedar Chief Scientific Officer and Richard Ho as CTO to round out the founding team deep technical and scientific expertise.

## **The a16z Investment Thesis**

Andreessen Horowitz leading a \$30M Series A extension in Protege in January 2025 signals strong conviction that data infrastructure will be foundational to AI advancement. The financing expanded the company's initial \$25M Series A from August 2024, bringing total funding to \$65M since founding in 2024. Returning investors include Footwork, CRV, Bloomberg Beta, Flex Capital, and Shaper Capital.

The thesis breaks down into several components. First, data access is genuinely a limiting factor for AI advancement right now. a16z's portfolio companies across AI and machine learning are all running into the same problem. They need diverse, high quality training data and cannot get it efficiently. Startups are burning millions in business development to cobble together datasets. Even well-funded companies struggle to access the data they need at the speed AI development requires. This creates demand for infrastructure that solves the problem systematically.

Second, the market is massive and growing. AI is eating every industry, and every application needs training data specific to its domain. Healthcare AI needs patient

data. Autonomous vehicles need driving data. Robotics needs sensor data from physical environments. Media companies need content libraries. The total market for training data could be larger than cloud computing because it cuts across every case.

Third, network effects create defensibility. Once Protege has relationships with hundreds of data suppliers and dozens of major AI companies, new entrants face enormous barriers. Data suppliers will not want to manage relationships with multiple platforms. AI builders will not want to integrate with multiple data sources when one platform gives them everything. The winner in this market could be a take-most, similar to how Snowflake dominated cloud data warehousing or how Databricks dominated data lakehouse architecture.

Fourth, the team has done this before. May built LiveRamp into a public company and Datavant into a \$7B business. Samuels brings operational expertise from leadership roles at both companies, including managing Datavant's Privacy Hub which was central to the company's compliance infrastructure. They know how to build new infrastructure platforms that scale across an ecosystem. They understand compliance, privacy, and trust at a deep level. They have relationships with data suppliers and know how to structure partnerships. This is not a team figuring it out for the first time. They are executing a proven playbook in a new market.

Fifth, timing is critical and favorable right now. AI companies are desperate for training data as public datasets run out. Frontier labs are willing to pay substantial amounts for unique datasets. Data suppliers are waking up to the value of their assets and looking for ways to monetize them. Regulatory frameworks around AI training data are still forming, creating an opportunity to help shape norms and standards. The window to build dominant data infrastructure is open but will not open forever.

The investment came from a16z's Bio and Health team, with partners Daisy Wolinsky and Eva Steinman involved. This makes sense given Protege's initial focus on healthcare data, though the platform has expanded into video, audio, and motion capture. T

Bio and Health team's involvement suggests a16z sees healthcare as the beachhead market but understands the platform will expand across verticals.

The \$30M round size on top of a previous \$25M suggests a16z expects Protege to grow quickly. This is not a seed investment in an unproven team testing product-market fit. It is a bet that the team can rapidly build supply and demand network effects before competitors emerge. The capital likely goes toward hiring engineers to build technical infrastructure, business development to sign data suppliers, sales to land AI consumers, and compliance infrastructure to operate across jurisdictions.

## What Protege Actually Does

Protege operates as a two-sided marketplace connecting data suppliers with AI builders, but calling it a marketplace undersells the technical and operational complexity involved. On the supply side, Protege partners with hospitals, health systems, labs, imaging centers, research networks, media companies, and other content holders. According to the company's announcements, Protege expanded its data partner network to hundreds of organizations in 2025, providing aggregated access to new data sources and formats.

Each partnership involves negotiating data licensing terms, building technical integrations to extract and normalize data, implementing privacy and compliance controls, and establishing revenue sharing arrangements. Protege provides revenue share payouts to data partners with each use, creating an economic incentive for content holders to contribute to the platform.

For healthcare specifically, Protege securely obtains patient data from multiple sources and stitches it into longitudinal, multimodal, anonymized patient-level datasets. This requires sophisticated entity resolution to match patient records across facilities without using identifiable information. A patient might have records at different hospitals, two labs, and an imaging center, all under slightly different names or spellings or with different identifiers. Protege's algorithms match these records probabilistically while maintaining HIPAA compliance through tokenization and other privacy-preserving techniques.

The data itself comes in wildly different formats. EHR data arrives as HL7 messages, FHIR resources, or proprietary formats depending on the source system. Lab results use LOINC codes. Diagnoses use ICD-10. Medications use RxNorm. Imaging data lives in DICOM files. Clinical notes are unstructured text. Protege normalizes all this into consistent schemas and data models that AI companies can actually use for training without building custom parsers for every data source.

Quality control happens at multiple stages. Protege validates data completeness, checks for anomalies, scores data quality, and flags potential issues before delivering datasets to customers. Bad training data causes model failures that might not surface until production, so quality assurance cannot be an afterthought. The platform tracks data lineage, versions datasets, and maintains audit trails for compliance purposes.

On the demand side, Protege serves frontier AI labs, AI application companies, and enterprises building internal AI capabilities. According to a16z's announcement, Protege already works with the majority of MAG7 public companies plus many private AI players. These companies use Protege to access curated datasets across healthcare, video, audio, motion capture, and other modalities without needing to negotiate hundreds of individual data partnerships.

The platform delivers data through multiple mechanisms depending on customer needs. Protege curates datasets from across its partner network to meet AI development needs, providing AI-ready data that integrates with modern ML workflows. The key value proposition is enabling AI builders to iterate quickly on model development rather than spending months or years on data acquisition and cleaning.

Beyond healthcare, Protege has expanded into other data modalities where similar problems exist. Media companies have vast archives of video and audio content that is valuable for training multimodal AI models but difficult to license at scale. Motion capture data from sports, entertainment, and research applications can train robotic and embodied AI systems. The same platform architecture that aggregates healthcare data can aggregate content libraries, with appropriate adjustments for different licensing and compliance requirements.

# Why This Team Can Execute

The reason a16z bet on this specific team rather than waiting for alternatives is that May, Samuels, and co-founders have unfair advantages that are nearly impossible to replicate. The relationship network alone is worth years of business development. May spent six years at Datavant building partnerships with major health systems, networks, imaging center chains, and health data vendors. Samuels spent over three years at Datavant in roles including General Manager of Privacy Hub, giving him deep relationships and operational knowledge. Before that, both had experience at LiveRamp building data partnerships in a different vertical.

Protege gets to inherit much of that relationship capital. When May calls the CEO of a major health system or lab network to discuss contributing data to AI training, he is not a stranger pitching a new idea. He is someone they worked with successfully at Datavant. The trust is already established. The credibility is already there. The understanding of their business constraints and compliance requirements is already built. This accelerates partnership development by an order of magnitude compared to a new team starting from scratch.

The technical expertise is equally important. Building infrastructure that can ingest data from thousands of heterogeneous sources, normalize it, maintain quality, preserve privacy, and deliver it at scale is genuinely hard. Most startups would spend years just building the technical foundation before delivering value to customers. Datavant's playbook, which May developed and Samuels helped execute, provides a blueprint for much of this. The core algorithms for entity resolution, the privacy-preserving techniques for tokenization, the data normalization pipelines, and the compliance infrastructure all transfer to Protege with modifications.

The regulatory and compliance knowledge is maybe the most underrated advantage. Healthcare data is among the most heavily regulated in the world. HIPAA has complex requirements around de-identification, business associate agreements, breach notification, and audit trails. Different states have additional privacy laws. International markets have GDPR and other frameworks. May and Samuels have years working with healthcare lawyers, privacy officers, compliance teams, and

regulators. They know what is permissible, what requires special handling, and how to structure agreements that satisfy all parties.

This expertise extends to understanding how different data holders think about sharing. Academic medical centers have research missions and IRB processes. Community hospitals care about revenue and liability. Labs want to protect competitive advantages. Health systems worry about patient trust. Protege can structure partnerships that address each institution's specific concerns because the team has done it hundreds of times before at Datavant.

The go-to-market approach also benefits from proven patterns. Datavant succeeded in becoming neutral infrastructure that everyone could use without creating competitive disadvantages. Health systems shared data through Datavant because it did not favor any single pharma company or technology platform. The same neutrality principle applies at Protege. AI companies will share a data platform if it does not give preferential treatment to competitors. Data suppliers will contribute to a platform that treats all customers fairly.

Operationally, the team knows how to scale quickly while maintaining quality. Datavant grew from zero to \$7B valuation in under four years through its merger with Ciox in 2021. The playbook involves starting with a wedge use case that delivers immediate value, proving the platform works, then expanding horizontally into adjacent datasets and vertically into new customer segments. Healthcare was the wedge for Protege. Video, audio, and motion capture are the horizontal expansion. Foundation model labs vs application companies vs enterprises represent vertical segmentation.

The engineering talent required to build this platform is also easier to recruit when the founders have successful exits and track records. Top data engineers want to work on hard problems with teams that have proven they can execute. Protege can attract senior technical talent from companies like Databricks, Snowflake, and Palantir by offering equity in a rocket ship with experienced founders who have built infrastructure companies before.

# Economic Realignment

The emergence of Protege and similar data infrastructure platforms shifts economic value throughout the AI stack in ways that are still playing out. For data suppliers, it creates new revenue streams that never existed before. Hospitals and health systems have always viewed patient data as a compliance burden and liability, not an asset. EHR systems cost millions to maintain, data teams prevent breaches, and sharing data opens up risk. But if you can monetize anonymized data for AI training while maintaining full compliance, suddenly that liability becomes valuable.

For healthcare providers specifically, the economics are compelling. A mid-size hospital system sitting on ten years of EHR data, imaging, and lab results represents significant value for training diagnostic models or clinical decision support systems. Previously, accessing that value required building internal data science teams, negotiating one-off partnerships, or simply leaving money on the table. Platforms like Protege that handle acquisition, anonymization, and licensing let providers generate revenue without adding headcount or compliance risk.

The revenue potential is meaningful relative to hospital margins. Health systems operate on thin margins, often 2 to 3% for non-profit hospitals. Adding a new revenue stream from data licensing, even if modest, can impact financial performance meaningfully. For struggling rural hospitals or safety-net providers, this could be the difference between staying open and closing.

Research networks and registries face similar dynamics. Organizations that collect patient outcomes data for specific conditions or treatments have spent years building these datasets for academic research. Now they can make that data available for commercial development with appropriate protections, creating funding that makes their core research mission more sustainable. Disease-specific registries, tumor boards, and clinical trial networks all sit on valuable longitudinal outcome data that AI companies desperately need.

Media companies and content owners are waking up to similar opportunities. Movie studios and broadcasters have massive video and audio archives that were previously

just sitting in vaults or used for limited internal purposes. Training multimodal models on diverse video content has enormous value for companies building conversation, video generation, or embodied AI systems. Licensing historical content for training creates a new revenue stream from otherwise dormant assets.

For AI builders, the economics flip from a major cost and bottleneck to a predictable expense. Instead of hiring business development teams to negotiate dozens of healthcare partnerships, burning six to twelve months on each, companies can access curated datasets through Protege in weeks. Instead of building internal data engineering teams to clean and integrate heterogeneous sources, they get normalized data ready for training. The time and cost savings are substantial, but the strategic value is larger.

Being able to iterate quickly on model hypotheses changes product development fundamentally. If you think adding a specific type of imaging data will improve diagnostic accuracy, you can test that in weeks rather than months or years. If a model performs poorly on certain patient populations, you can quickly source additional training data to address the gap. Speed of iteration becomes a competitive advantage and Protege enables that speed.

Pricing models will be critical for how this plays out. Traditional enterprise data deals involve lengthy negotiations, volume commitments, and opaque pricing. That works for established companies with data budgets but kills startup experimentation. Protege can offer transparent, usage-based pricing aligned to startup economics, which enables a much broader set of AI builders to access valuable training data. This is similar to how AWS democratized infrastructure access compared to buying your own servers.

There are interesting dynamics around data exclusivity and competitive advantage. Should leading AI companies be able to license exclusive access to certain datasets? Does that create unfair advantages, or is it just normal competitive tactics? Protege needs to balance enabling competition with allowing differentiation. The likely equilibrium involves a mix of widely available datasets that level the playing field.

exclusive arrangements for unique data sources, similar to how cloud infrastructure works today.

The revenue split between Protege and data suppliers also matters. If Protege takes too much margin, data suppliers will try to go direct or use competing platforms. If Protege gives away too much margin, the business will not be sustainable or profitable enough to justify a16z's valuation expectations. The right split probably varies by data type, exclusivity, and supplier bargaining power. Large health systems have more leverage than small research networks. Unique datasets command better economics than commoditized data.

## What This Means for Builders

For founders building AI companies, the implications of mature data infrastructure shift strategic priorities in several ways. Data strategy moves from being primarily a business development and operations challenge to being a product and engineering question. Instead of hiring salespeople to negotiate hospital partnerships, you hire ML engineers to evaluate dataset quality and design training pipelines. Instead of building custom ETL for each data source, you integrate with standardized APIs.

This lowers barriers to entry for new AI applications that were previously too difficult for startups to pursue. Building a diagnostic radiology model used to require years of hospital partnerships before training the first model. Now you can get started in weeks. This opens up entire categories of healthcare AI that were only accessible to well-funded, experienced teams. The same pattern will play out in other verticals as Protege and similar platforms expand beyond healthcare.

Competitive dynamics shift toward model architecture, training techniques, and application-specific optimization rather than pure data access. When everyone can access similar baseline datasets, differentiation comes from what you do with the data. This is probably healthier for innovation overall, since it rewards technical capabilities rather than just partnership skills. Companies compete on actual AI capabilities instead of who negotiated better data deals.

For investors, data infrastructure platforms represent a different risk-return profile than typical SaaS businesses. Network effects are strong once you have critical mass on both supply and demand sides. Marginal costs for incremental data sources and customers are relatively low compared to initial platform development. Switching costs are moderate to high once AI companies integrate data pipelines into training workflows. The business model looks more like marketplace economics than traditional software.

Revenue concentration around a small number of large AI customers creates risk but also validates product-market fit. If the leading frontier model builders all use your platform, that proves the value proposition strongly. The question becomes whether you can expand beyond anchor customers to serve the long tail of AI builders. Palantir is already working with majority of MAG7 companies plus large private AI players, suggesting they have the anchor customers locked in. Expanding to mid-market and smaller companies will determine ultimate market size.

Regulatory risk is material but manageable with proper compliance infrastructure. Healthcare data in particular faces ongoing regulatory scrutiny, and rules around training data are still evolving. Platforms with deep compliance expertise and relationships with regulators have advantages, but everyone operating in this space accepts some regulatory uncertainty. The team's Datavant experience navigating HIPAA, state privacy laws, and international frameworks is a significant derisking factor for investors.

International expansion creates obvious growth opportunities but requires rebuilding supply-side relationships and compliance infrastructure in each major market. Different localization requirements, different privacy regimes, and varying healthcare system structures mean you cannot just flip a switch to operate in Europe or Asia. Companies that can execute international expansion effectively will build significant moats. Amazon and Samuels did this at Datavant, giving them playbooks for how to approach different markets.

The broader pattern extends beyond healthcare to any domain with valuable private data. Manufacturing and industrial companies with sensor data from physical

processes could enable embodied AI and robotics. Financial institutions with transaction data could train better fraud detection and risk models.

Telecommunications companies with network data could improve infrastructure optimization. The playbook established in healthcare likely applies across multiple verticals, each of which could be as large as healthcare alone.

What remains uncertain is whether data infrastructure becomes a winner-take-all market or supports multiple specialized platforms. Arguments exist on both sides. Network effects and economies of scale in building supply relationships favor concentration. But vertical specialization, regional focus, and different data modalities might support multiple winners. Healthcare alone might sustain several platforms focusing on different data types or customer segments. The next few years will determine market structure.

The timing question matters significantly. Data infrastructure platforms that establish themselves now, while frontier AI labs are desperate for training data, will be sticky even as the market matures. Companies that wait risk entering a market with established incumbents and locked-up supply relationships. For entrepreneurs, the window to build in this category is open but probably measured in quarters, not years. Protege's \$65M in funding and a16z backing will accelerate their timeline and make it harder for followers to catch up.

Ultimately, the shift from compute and architecture to data as the primary constraint in AI advancement represents a fundamental market transition. The companies that build infrastructure to make real-world data accessible at scale are not just service providers but enablers of the entire next wave of AI applications. That is rare positioning and, if executed well, creates durable value. Travis May has built this kind of infrastructure company twice before with LiveRamp and Datavant. The a16z bet is that he and Bobby Samuels can do it again in a market that is potentially larger than their previous two companies combined.



4 Likes

[← Previous](#)

[Next](#)

## Discussion about this post

[Comments](#)

[Restacks](#)



Write a comment...