

The Trust Factor: Delivering Verifiable & Confidential Intelligence

EXECUTIVE SUMMARY

AI is no longer a peripheral productivity tool; it is becoming infrastructure. As models gain memory, autonomy, and influence over decisions, intelligence is compounding faster than trust. Enterprises are asked to surrender their most valuable asset (proprietary data) simply to access insight, often without any verifiable guarantee of how that data is handled. Every document uploaded and every prompt sent, becomes an irrevocable act of trust.

Centralized AI systems operate on assertions rather than proof. Users are expected to believe that data is private, that sensitive inputs are not retained or used for training, and that access controls are enforced correctly across opaque infrastructure. These assurances are contractual and policy-based, not cryptographic. They cannot be independently audited or verified, and they break down when incentives shift or stakes rise.

This case study explores an alternative: a decentralized AI-native assistant that replaces custodial trust with verifiable cryptographic guarantees. By enforcing integrity across compute, context, and collaboration, intelligence can be delivered online without surrendering data ownership. Models can reason over sensitive information while neither model operators nor infrastructure providers can access it. Trust is no longer assumed, it is enforced.

TABLE OF CONTENTS

- [1. WHY TODAY'S CENTRALIZED AI FAILS AT TRUST](#)
- [2. PRODUCT CONCEPT](#)
- [3. WHY THIS WINS](#)
- [4. VC APPLICATION / PROOF OF CONCEPT](#)
- [5. CONFIDENTIAL COMPUTE - GPU TEE](#)
- [6. DATA STORAGE](#)
- [7. PRODUCT 1 - CONFIDENTIAL "CHATGPT"](#)
- [8. PRODUCT 2 - CONFIDENTIAL CUSTOM "GPT's"](#)
- [9. FUND DATA RETRIEVAL AUGMENTED GENERATION \(RAG\)](#)
- [10. DECENTRALIZED COMPONENTS DEEP DIVE](#)
- [11. KEY TAKEAWAYS](#)

1. WHY TODAY'S CENTRALIZED AI FAILS AT TRUST

Macro Shift: AI Is Becoming Infrastructure

Modern AI systems force organizations into a false trade-off between capability, usability and trust. Centralized platforms deliver polished interfaces, fast inference, and powerful reasoning, but they do so by requiring full custodianship over user data. Every prompt, document, and contextual signal is processed on infrastructure the user cannot inspect, audit, or control. Privacy assurances are enforced through policies and terms of service, not cryptography. Enterprises are asked to trust that data is not logged, retained, reused, or accessed improperly, yet none of these guarantees are independently verifiable. As AI becomes embedded in strategic and operational decision-making, this model of asserted trust becomes structurally inadequate.

Open Source AI: Verifiability Without Usability

Open-source large language models were expected to resolve this tension. In theory, they offer inspection, self-hosting, and full data control. In practice, adoption has stalled. Running models locally introduces real operational costs: constrained hardware, degraded performance, maintenance overhead, and rising latency as context grows. Open-source deployments lack the application layer that made centralized AI dominant and as a result, deployments are often stuck with smaller models, weaker reasoning and higher hallucination rates. The limitation is not philosophical; it is experiential. Faced with slower responses and weaker usability, organizations predictably optimize for output quality and convenience over theoretical privacy gains.

The Missing Option

The result is a gap that neither centralized nor local approaches can fully solve. Most enterprises already operate in the cloud, with distributed teams, shared context, and SaaS-native workflows. Fully local AI stacks introduce friction, duplication, and operational complexity, while failing to support collaboration or scale. Centralized AI, meanwhile, offers excellent usability but demands unconditional trust in opaque operators. What organizations actually want is a cloud-like AI experience: high performance, rich context, multi-model access and collaboration - without surrendering data or relying on unverifiable promises.

Decentralized AI-native systems make this possible by relocating trust away from operators and into verifiable hardware and cryptographic enforcement. By combining confidential compute, GPU-based TEEs, sovereign user-owned memory, and verifiable data layers, intelligence can be delivered online without custodial risk. Models can reason over sensitive data without exposing it to hosts or infrastructure providers. Context remains portable, user-owned, and revocable. Every execution can be independently verified, end-to-end. This is not decentralization for ideology's sake. It is an operational response to a simple reality: AI is becoming agentic and embedded in high-stakes workflows. In that environment, privacy cannot be a promise, execution cannot be assumed, and trust cannot be delegated. Verifiable privacy and data sovereignty are no longer trade-offs. They are prerequisites.

2. PRODUCT CONCEPT

This product exists because policy-based trust is no longer sufficient. Modern AI has reached a point where its capabilities rival human judgment in high-stakes environments, yet the trust model underpinning it remains centralized, opaque, and permissioned. Users are asked to surrender ownership of their data in exchange for convenience, governed by contracts they cannot verify and systems they cannot inspect.

We reject that trade-off.

This is a cloud-native AI assistant that delivers the usability, performance, and collaboration people expect from modern centralized AI, without custodial trust. Instead of policies, promises, or access controls, trust is enforced at the hardware level, and proven cryptographically. Every interaction is verifiable. Every guarantee is provable.

The user does not trust the platform.
The platform proves itself to the user.

Product 1: Confidential “ChatGPT”

Stack: Confidential Compute + Confidential GPU within Trusted Execution Environments (TEE)

Confidential ChatGPT is a general-purpose AI assistant comparable to ChatGPT or Claude but built on an entirely different philosophy. Prompts, context, and model execution remain encrypted at rest and in use. No operator, host, cloud provider, or infrastructure intermediary can access user data.

This is not privacy as a setting. This is privacy enforced by silicon.

Every inference produces a cryptographic proof that the model executed inside a genuine TEE. Users can independently verify that memory was isolated, execution was correct, and no third party had access to context during runtime. Trust is no longer implied, it is mathematically demonstrated.

This removes the false dichotomy between usability and privacy. Users retain the speed, quality, and ergonomics of modern cloud AI while gaining guarantees that their data was never exposed, logged, or retained. Confidential ChatGPT is designed for environments where leakage is unacceptable: regulated teams, internal strategy, financial analysis, legal review, and sensitive research. Privacy here is not a policy promise, instead, is a property of the system.

Product 2: Confidential “Custom GPT’s”

Stack: Confidential Compute + Confidential GPU TEE + Sovereign Context + Data Lake / Warehouse

Confidential Custom GPTs extend these guarantees to configurable, shareable AI agents. Instructions, long-term memory, and private datasets can be combined into purpose-built assistants, without ever centralizing ownership or control.

This directly challenges existing Custom GPT frameworks, where uploaded documents, context, and memory are permanently custodial. In this model, context remains sovereign. Access can be granted, shared, or revoked without exposing raw data. Every query and response remains cryptographically verifiable end-to-end.

Confidential Custom GPTs occupy the missing middle ground between centralized AI and fully local systems. They preserve persistence, sharing, and usability - while enforcing cryptographic guarantees over execution, memory, and data access. This makes them suitable for internal research assistants, diligence copilots, LP-facing Q&A agents, and collaborative workflows where trust assumptions must be minimized.

This is not incremental security layered on top of centralized AI. It is a fundamentally different trust architecture - one where intelligence can scale without central custody.

3. WHY THIS WINS

Compared to centralized AI platforms, the system provides no custodial access to prompts, documents, or memory. It does not depend on policy assurances or operator integrity. Execution, privacy, and correctness are independently verifiable.

Compared to open-source or fully local LLM deployments, it removes the burden of hardware management, performance degradation, and ongoing maintenance, while preserving native collaboration, persistent context, and multi-model access.

Compared to so-called “secure AI wrappers,” this is not encryption layered onto centralized infrastructure. Trust is enforced at execution, not abstracted at the application layer.

Rather than relying on a single vendor or centralized operator, the system is modular by design. It integrates leading confidential providers into a unified trust model, avoiding vendor lock-in and eliminating single points of failure.

Companies featured in this case study include Phala, Space and Time, BlueNexus, AO, NEAR, Chutes, Tinfoil, Memvid, XTrace, Redpill, and Irys.

The result is military-grade protection for every AI interaction without sacrificing usability, collaboration, or performance. Intelligence is delivered online, at scale, with the ergonomics of modern cloud AI and the guarantees of sovereign execution. This is not an additional security layer on top of centralized AI, it is a fundamentally different trust architecture.

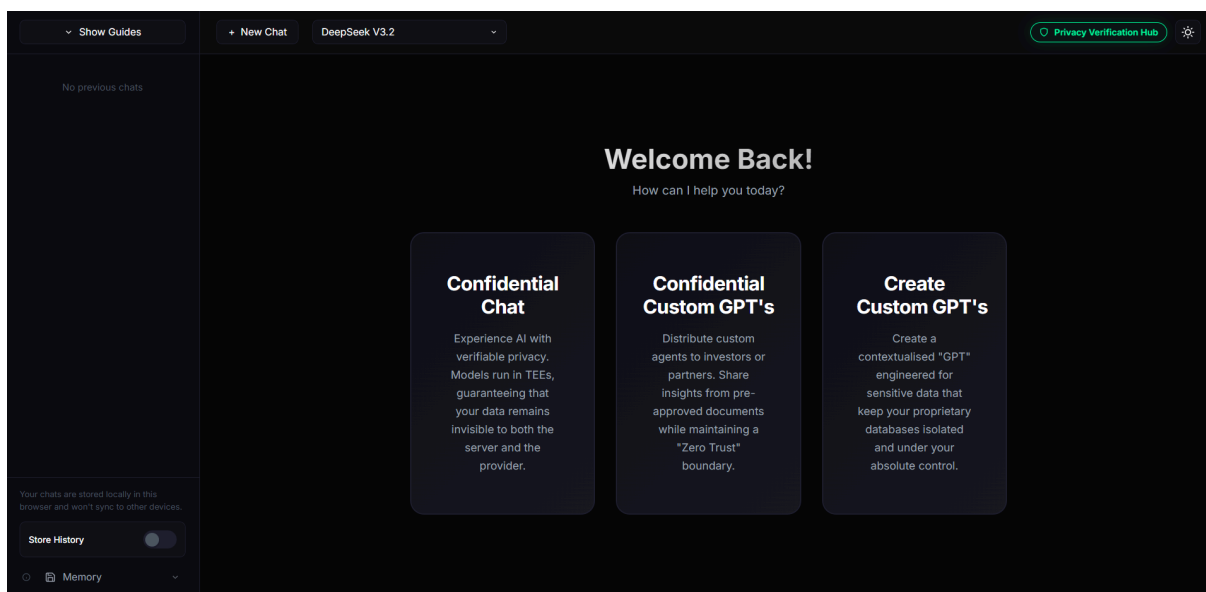
This is the trust architecture required for AI at institutional scale.

4. VC APPLICATION / PROOF OF CONCEPT

This use case was designed for operators within financial services, specifically venture capital, private equity, and hedge funds that actively raise capital. The product supports multiple core workflows. First, it enables teams to upload, query, and reason over investment documents to support diligence, analysis, and internal decision-making. Second, it allows firms to share pitch decks and related materials in a secure, controlled environment without compromising confidentiality.

As this is an internal-facing application, several layers of the broader decentralized AI stack can be intentionally omitted. At launch, the application presents a unified interface exposing the core product concepts. Custom Confidential GPTs are preconfigured, saved, and persistently available, allowing team members to access and share them on demand without repeated setup. This demonstrates that confidential AI is not theoretical - it can be deployed today to support real institutional workflows without reintroducing custodial risk.

Homescreen / Product Landing Page



Each time “+ New Chat” is selected, the user is returned to the home screen. From there, activities can be selected, models switched, privacy guarantees, memory managed, and all guides are accessible.

5. CONFIDENTIAL COMPUTE - GPU TEE

Model Availability

For model selection, only confidential AI models running inside GPU TEEs are available, sourced from multiple ecosystems (Phala, Tinfoil, Near, and Chutes). For a deeper explanation of TEEs, we recommend reading this [article](#).

DeepSeek V3.2

PHALA NETWORK		TINFOIL	
all-MiniLM-L6-v2	GPU TEE	DeepSeek R1 0528	GPU TEE
DeepSeek V3.2	GPU TEE	Kimi K2 Thinking	GPU TEE
Gemma 3 27B	GPU TEE	Qwen3 Coder 480B A35B	GPU TEE
GPT OSS 120B	GPU TEE	Qwen3 VL 30B A3B Instruct	GPU TEE
GPT OSS 20B	GPU TEE	NEAR AI	
Llama 3.3 70B Instruct	GPU TEE	DeepSeek V3.1	GPU TEE
Qwen2.5 7B Instruct	GPU TEE	Qwen3 30B A3B Instruct	GPU TEE
Qwen2.5 VL 72B Instruct	GPU TEE	Z-AI GLM 4.6	GPU TEE
Qwen3 Embedding 8B	GPU TEE	CHUTES	
Venice Uncensored 24B	GPU TEE	MiniMax M2.1	GPU TEE

Models are grouped by their respective providers, while GPU TEE hosting is delivered via API by Redpill, a product of [Hashforest](#), which also operates Phala. Redpill provides verification via API across the following components:

What Each Check Means

Check	What It Proves
Intel TDX quote verified	Code runs in genuine Intel TDX CPU enclave
Report data binds signing address	Signing key is generated inside TEE
Report data embeds request nonce	Attestation is fresh (not replayed)
GPU payload nonce matches	GPU attestation is for this specific request
NVIDIA attestation verdict	GPU is genuine H100/H200 with TEE
mr_config matches compose hash	Running code matches the Docker compose shown

Source: [Redpill](#)


Confidential VM


This, combined with a confidential server hosted on Phala Confidential Virtual Machine (CVM), is critical for enterprise and institutional use cases. Without GPU TEEs, confidential compute stops at the CPU boundary and breaks where AI workloads concentrate. Without confidential VMs, GPU TEEs become isolated islands surrounded by exposed logic. Together, confidential VMs and GPU TEEs enforce end-to-end confidentiality, integrity, and non-access, allowing models to reason over sensitive data online without reintroducing custodial trust.


To maximize trust, the Privacy Verification Hub is designed so that GPU and CPU attestations can be independently verified on every model change and every refresh.


[illegible]

In addition, Phala provides a TEE Attestation [Explorer](#), allowing users to copy the TDX quote, paste it into the explorer, and independently validate the attestation.



PHALA


AUTOMATA


zkVerify


VERIFIED

The attestation is verified and safe to use.


Proof of Cloud: This device is not in our verified facilities registry. [Learn more](#)

REPORT DATA

Show Hex

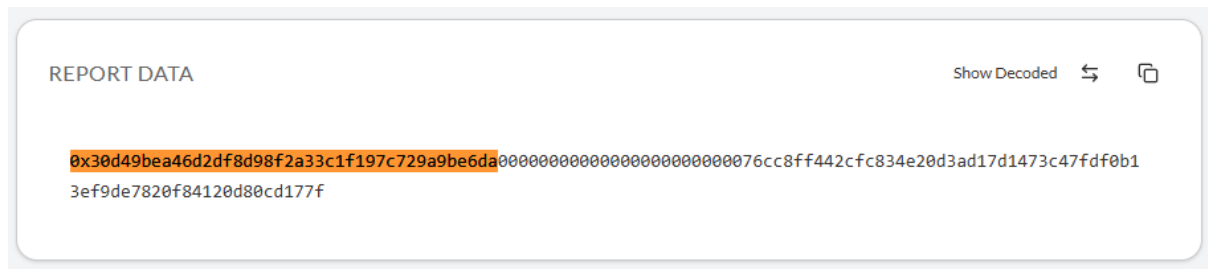
↔

📄

0qF^<|r

0B4:}Bs

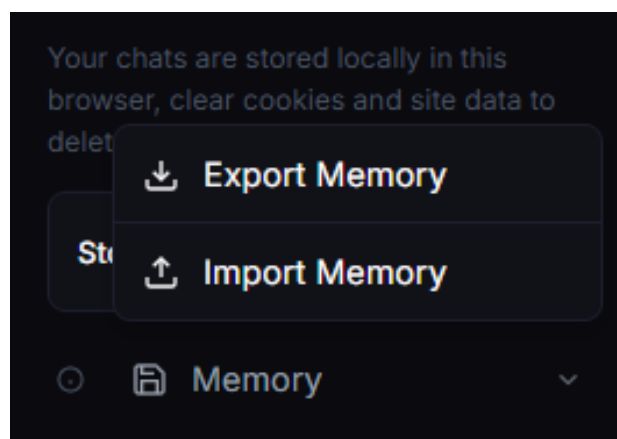
Selecting “Show Hex” decodes the report data. The signing address can then be copied from the Privacy Verification Center and verified directly against the decoded report.



6. DATA STORAGE

Portability

Data portability should not be a barrier. To ensure full user control, chats can be easily imported or exported in the same format used by ChatGPT. By default, conversations are stored locally in the browser; clearing cookies and site data will remove all saved chats. Users have multiple options for how their data is stored. They may choose to keep conversations stored locally on their device and periodically export them for archival or compliance purposes. However, when stored purely locally, these conversations are not available to the AI model for retrieval-augmented generation (RAG); they function only as historical records rather than active memory.



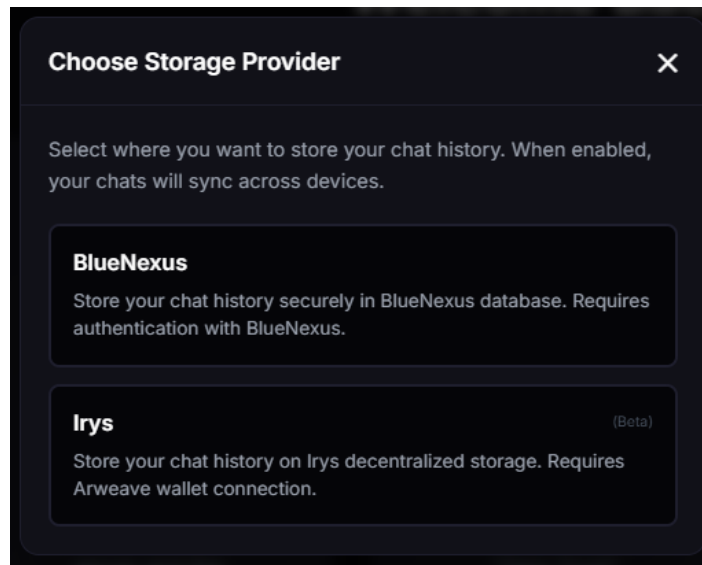
XTrace

For persistent, model-accessible memory, one option is the XTrace Personal Memory Vault. XTrace is a user-controlled, encrypted memory layer that allows individuals to selectively provide long-term context to their private AI models without surrendering custody of the underlying data. Through the XTrace [Chrome extension](#), conversations, documents, and contextual signals are securely stored and made available for RAG at inference time, while remaining fully owned and controlled by the user. The AI can reason over this memory, but neither the hosting infrastructure nor third-party operators can access or extract the data.

This approach enables continuity, personalization and collaboration across sessions, without reintroducing the custodial risk inherent in centralized AI memory systems.

In-Solution Storage Options

If the “Store History” option is selected, this prompt is displayed and the system automatically saves and uploads chat history to a BlueNexus or Irys account. Because this is an internal tool, these storage options are not tied to a default login flow; setup and access requirements are configured individually for each environment.



Blue Nexus

BlueNexus provides multiple ways for developers to integrate user connections into their platform. There are three primary integration models, each suited to different levels of control and user experience.

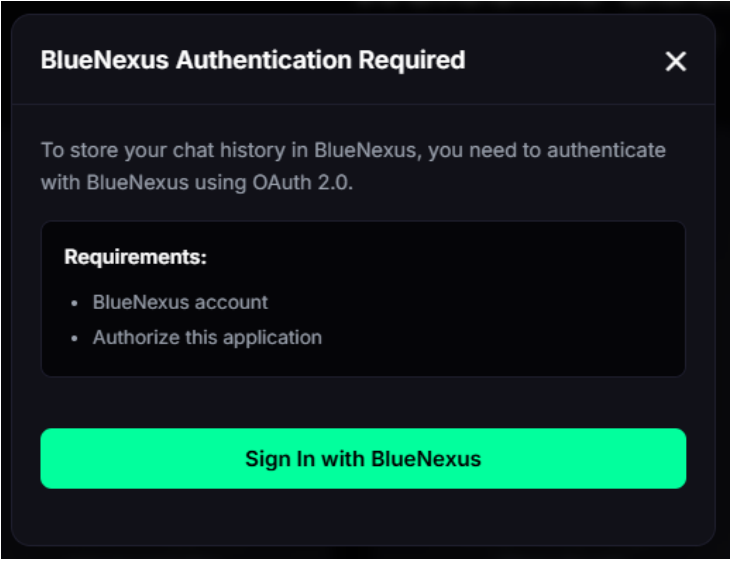
Sign in with BlueNexus - Use BlueNexus accounts as the authentication solution.
Connect BlueNexus - Allow existing users to connect their BlueNexus accounts.
White-Label Integration - Programmatically create and manage BlueNexus accounts.

This approach offers several benefits for both developers and users:

- Delegates user management, authentication, and identity storage to BlueNexus.
- Fastest setup for new applications or platforms without an existing auth system.
- Users maintain full control over their accounts and credentials.
- Users retain complete sovereignty and can revoke access at any time.

White-Label benefits:

- Users never interact with the BlueNexus login interface.
- Seamless, fully branded user experience.
- Fine-grained control over user data and account creation.
- Best suited for large platforms embedding BlueNexus as a backend service.



Further BlueNexus Integrations Opportunities:

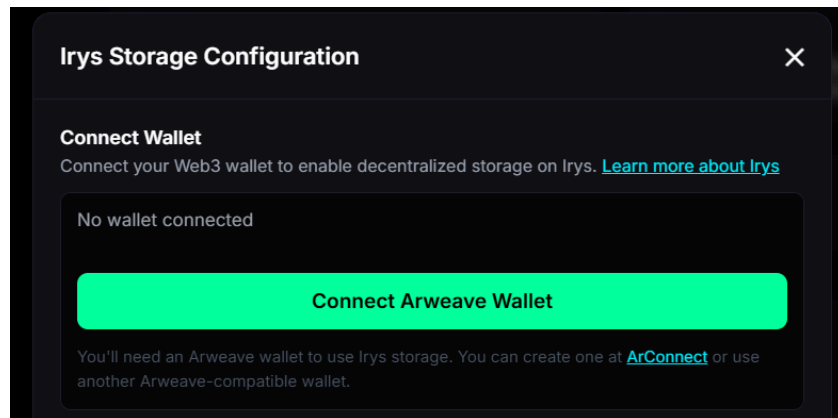
Access Method	When to Use	Key Characteristics
1. AI Connectivity (via MCP)	Connect your LLM to real-time user data. Ideal for chat-bots, assistants, and generative-AI workflows.	<ul style="list-style-type: none">• MCP Service for all third party data sources.• Automatic handling of authentication, scaling on managed infrastructure.
2. Third-Party Services	Pulling data from external services (Health platforms, CRMs, analytics platforms, SaaS tools, etc.) that a user has connected to their account.	<ul style="list-style-type: none">• Unified RESTful API layer that abstracts each external provider.• Automatic token handling, pagination, and rate-limit management.
3. Database Storage	Persistent, structured or unstructured data that is stored within a user account.	<ul style="list-style-type: none">• JSON-document store (MongoDB-compatible).• Unlimited <i>collection</i> per user• Data schema support for enforcing structured data (optional)• Shared across all apps that belong to the same user account.

Source: [BlueNexus](#)

Within product 1 (Confidential ChatGPT), this opens the door to integrating third-party services, allowing private models to securely query external data sources inside a confidential execution environment. For more details on BlueNexus security and privacy guarantees, refer to their [documentation](#).

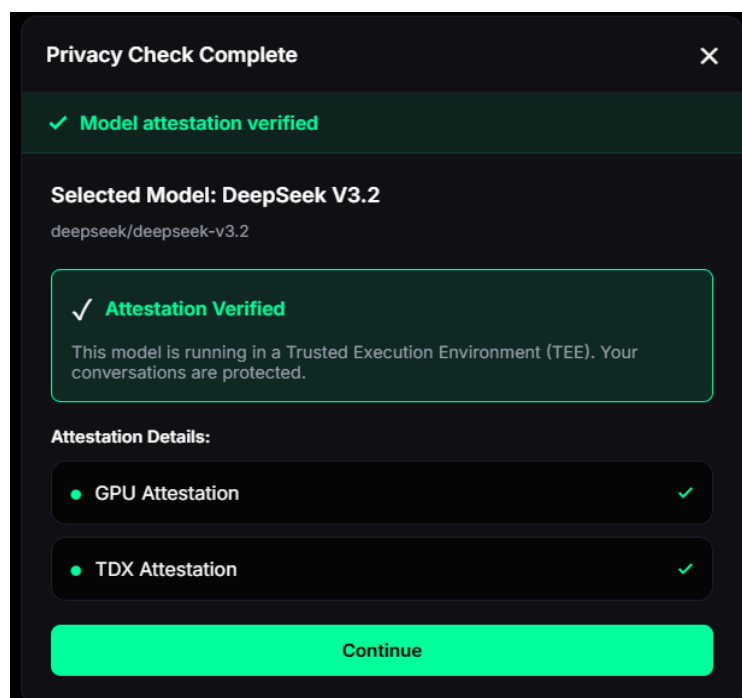
Irys (Arweave)

Irys can also be used as a storage option. It is not a blockchain itself, but a layer built on top of Arweave. Irys bundles files efficiently and stores them permanently on Arweave: think of it as a decentralized, permanent flash drive. While files can be uploaded directly to Arweave, the process is often expensive and complex. However, this storage is permanent so this is positioned as an optional feature for storing only the most important files, rather than all data by default. You can save just 1 chat and upload that to your knowledge.



7. PRODUCT 1 - CONFIDENTIAL "CHATGPT"


After selecting this option on the homescreen, users are prompted to confirm their model selection. A dropdown menu provides recommendations based on task suitability. A further privacy verification step confirms both GPU and CPU attestations before the session begins.



Four options become available:


How do you want to start this chat?

Choose how you want to provide context. Upload files, use a DocSend deck, load a Memvid knowledge base, or start chatting.




Upload files

Attach PDFs or docs and let the model use them as context.




Use DocSend link

Fetch a DocSend deck and use it as the knowledge base.



Upload Memvid Mv2 file

Load a Memvid knowledge base (.mv2) file and chat with its contents.




Just chat

Skip documents and start a normal confidential chat.

Upload Files / DocSend Flow:

- User uploads a PDF → encrypted and stored on server (Confidential VM - Phala).
- User submits a question → PDFs are decrypted in memory, searched and relevant chunks extracted.
- Context is sent to the TEE GPU model → model generates response using PDF.
- PDFs are automatically deleted after one hour.

 **1 PDF file uploaded**
✓ Encrypted before upload • ✓ Processed in GPU TEE • Never stored in plaintext • Temp storage (deleted within 1hr)

Chat with your documents

DocSend Link

DocSend Link

https://docsend.com/view/...

Email (Optional - for password-protected documents)

your@email.com

Password/Passcode (Optional - for password-protected documents)

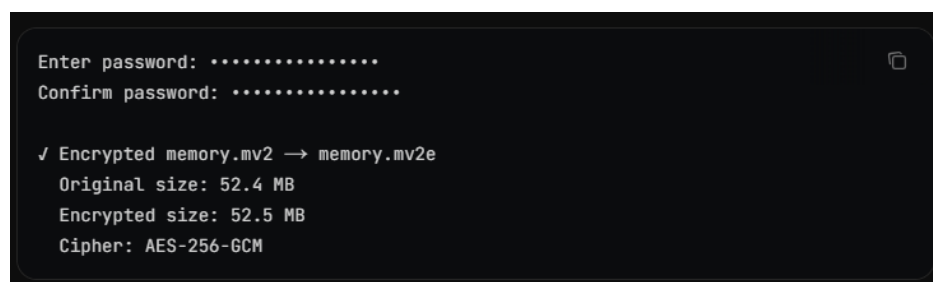
Enter password

Download from DocSend

For DocSend, users can input the link and optionally provide credentials if the document is protected. If the DocSend owner has disabled file downloads, the flow will not function. Otherwise, the process mirrors that of a standard PDF upload.

Memvid

Memvid converts AI memory into a single, portable, serverless file by embedding text and metadata directly into video frames. This eliminates the need for traditional RAG pipelines or vector databases, giving developers a fully portable memory layer. Users can upload documents within “Create Contextualized GPTs”, generate a .mv2 file, optionally encrypt it with AES-256-GCM, and download it. The file can then be shared, maintaining user ownership and allowing others to securely query its contents without exposing sensitive data.



Source: [Memvid](#)

Ingesting Mv2 File

To create a .mv2 file, select the “Share Contextualized Agent” option from the homescreen UI. Once the documents are uploaded and the file is generated, it can be reloaded into “Confidential ChatGPT”, making all embedded knowledge immediately queryable. This design ensures seamless portability: files can be copied, synced across devices, or used fully offline, delivering instant retrieval and persistent memory without requiring servers, databases, or complex infrastructure.

```
C:\Memvid>node ingest.js
=== Building High-Fidelity AI Memory ===
-> Processing PDF...
-> PDF Parser failed. Using expanded text fallback for testing.
-> Creating 24 overlapping memory frames...

Building search index...
--- SUCCESS: Full Memory Built ---

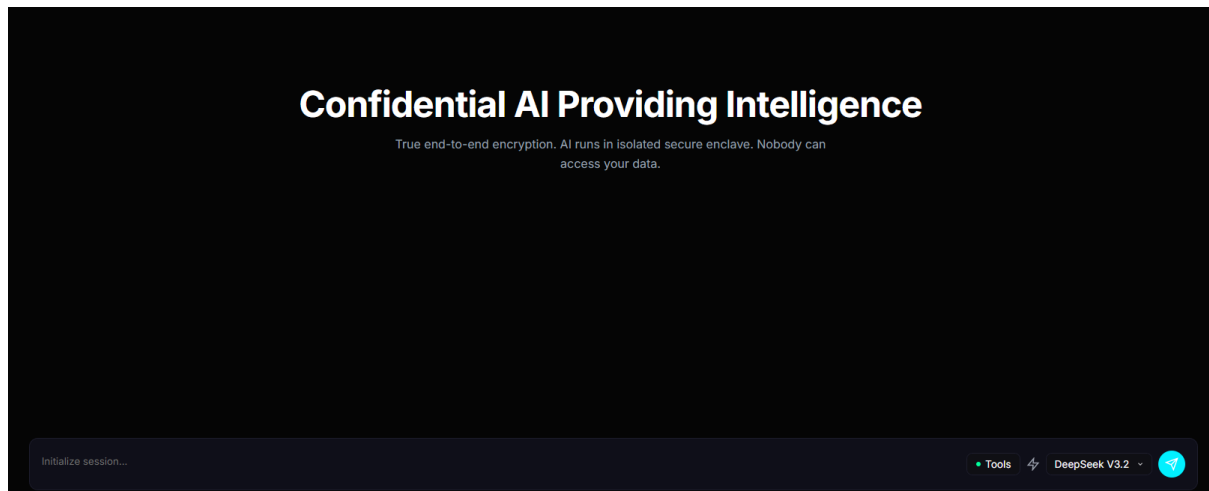
C:\Memvid>node fixsearch.js
Searching for allocation details...
✗ No matches found. Try running 'timeline.js' to see what's inside.

C:\Memvid>node timeline.js
=== Accessing Memvid Timeline ===
Total Frames Stored: Unknown

Displaying 24 Knowledge Frames:

[FRAME #1] -----
URI:      mv2://fund-d-q4/segment-0
CONTENT:  FUND D - SPECIAL SITUATIONS Q4 FACTSHEET.
STRATEGY: Focus on high-alpha opportunities in the distressed crypto ecosystem
CREATED:  06/01/2026, 21:27:55
[FRAME #2] -----
URI:      mv2://fund-d-q4/segment-1
CONTENT:  stem.
TARGET SIZE: $25,000,000. MINIMUM COMMITMENT: $100,000.
FEES: 2% Management Fee, 20% Performance Fee with High-Wat
CREATED:  06/01/2026, 21:27:55
[FRAME #3] -----
URI:      mv2://fund-d-q4/segment-2
CONTENT:  -Water Mark.
LOCKUP: 2-Year initial lockup period for all limited partners.
ALLOCATION: 60% Distressed Tokens, 30% Restr
```

After this we are now at the chat interface where you can now use AI with cryptographic guarantees. Rest assured, your data is never retained and never used for training.

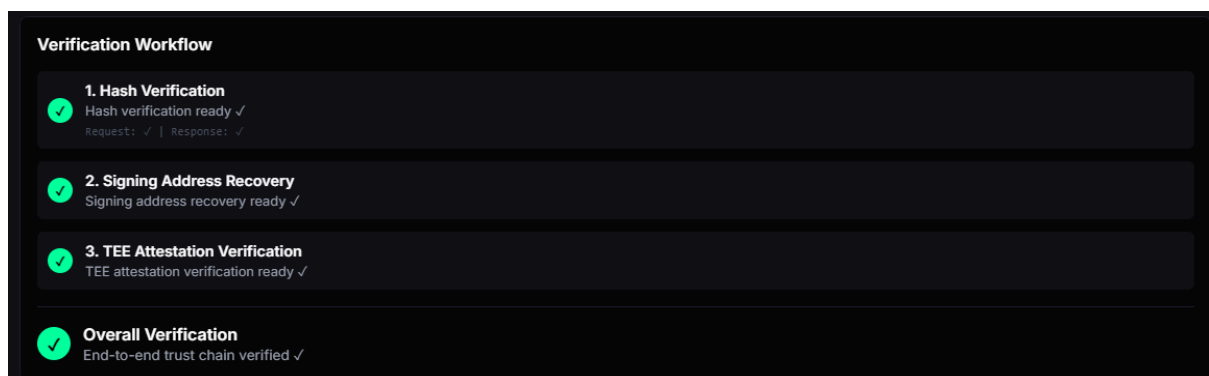


Every response can be copied to clipboard, regenerated with the same model and independently verified, including full verification or confirmation of the model itself.



Message Verification

Message verification provides additional details, such as cryptographic signatures, giving you complete confidence in the integrity of each interaction.

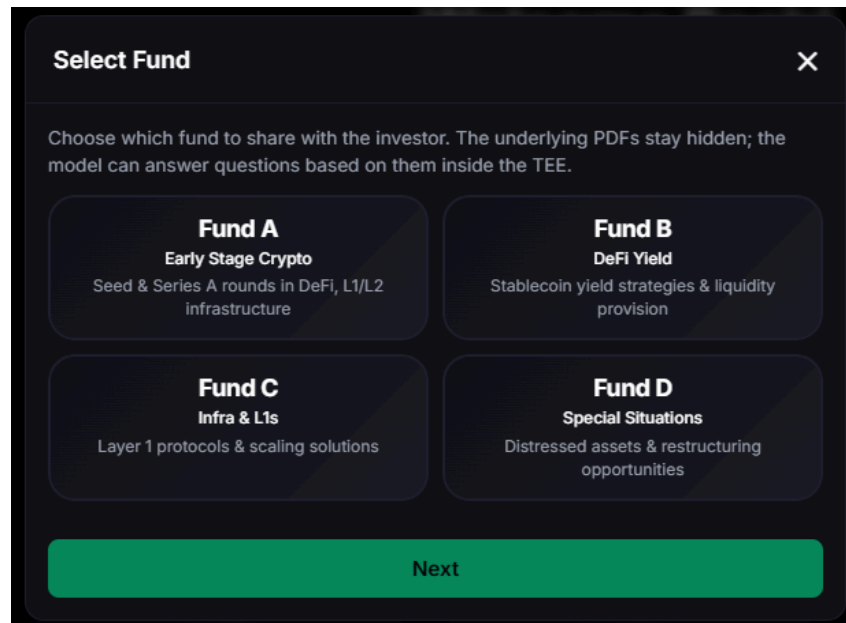


To fully verify a signature, you need to:

- Verify request/response hashes match.
- Recover the signing address from the signature.
- Fetch fresh attestation for that signing address.
- Verify the complete attestation (proving the signer is genuine TEE).

This creates a trust chain from your response → signature → signing key → TEE hardware.

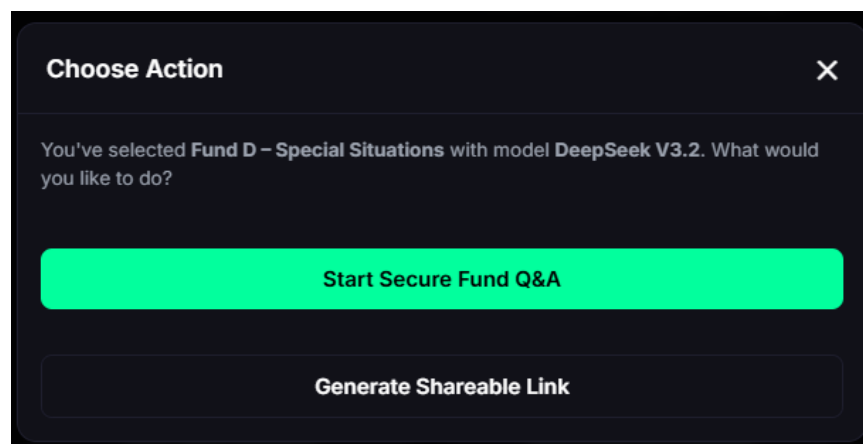
8. PRODUCT 2 - CONFIDENTIAL CUSTOM GPT's



For this case study, we set up four funds. This is the first stage after selecting the “Confidential Custom GPT’s” option. As this platform is designed for fund managers (with no coding experience), all they would have to do is access the code - which in this case is hosted on a confidential VM via Phala and update the PDFs or JSON files containing the relevant information.

Two options are now available:

1. Start Secure Fund Q&A: This allows administrators to quickly locate fund terms, ask questions on behalf of the original querier and provide a response.
2. Generate a Shareable Link: This creates a link that can be shared with an LP, enabling them to ask questions directly. This feature includes multiple customization settings. When you generate a shareable link, your chosen settings are encoded directly into the URL. If the link has expired, access is denied and an appropriate error message is displayed.



Additional Settings for Shareable Link

×

Configure what features will be available to users who access the shareable link.

[Learn how shareable links work and how expiration is enforced](#)

Allow them to see what pages

Users can view which pages/documents are available

☒ See Confidential AI Guide
 ☐ See Decentralised Fund Agent
 ☐ See Privacy Verification Hub Guide

☐ Create New Chat

⚠ They can have access to chat / create contextualised agent / share all fund agents

☐ Allow light mode

Users can switch between light and dark themes

How long will the link be valid

Set an expiration time for the shareable link

☐ Enable

☒ Allow export memory

Users can export their conversation history

☐ Allow import memory

Users can import conversation history









☒ Allow them to change models

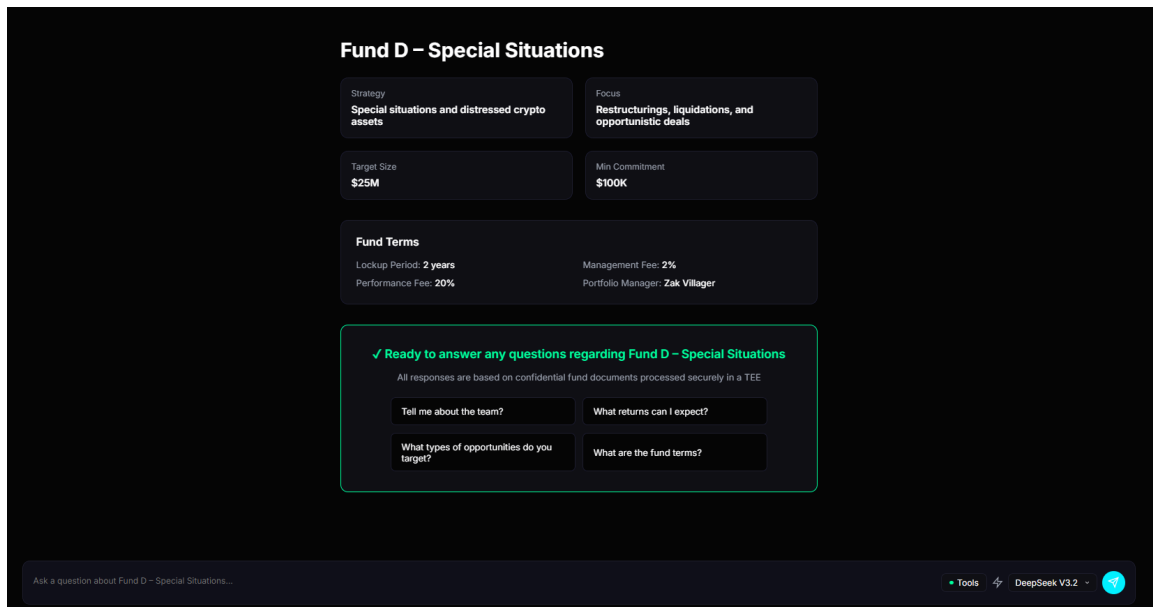
Users can select different AI models for the conversation

Generate Shareable Link

Update Fund Data

The codebase is organized around four fund folders - one for each fund. To update or add information, place the relevant documents and files into the corresponding folder.

 fundA	22/12/2025 21:17	File folder
 fundB	22/12/2025 21:17	File folder
 fundC	22/12/2025 21:17	File folder
 fundD	22/12/2025 21:17	File folder
 fundA	23/12/2025 20:40	JSON File
 fundB	23/12/2025 20:40	JSON File
 fundC	23/12/2025 20:40	JSON File
 fundD	23/12/2025 20:40	JSON File



To update any data on the fund UI, simply navigate to the corresponding folder and modify the relevant .json files. This allows you to quickly adjust fund details or update the questions based on the most common or recent queries, making it fast and easy to keep the information accurate and accessible.

```
{
  "id": "fundD",
  "name": "Fund D – Special Situations",
  "description": "Special situations / distressed opportunities with supporting docs.",
  "information": {
    "strategy": "Special situations and distressed crypto assets",
    "focus": "Restructurings, liquidations, and opportunistic deals",
    "targetSize": "$25M",
    "minCommitment": "$100K",
    "lockupPeriod": "2 years",
    "managementFee": "2%",
    "performanceFee": "20%",
    "team": "Zak Villager",
    "keyHighlights": [
      "Opportunistic approach to distressed assets",
      "Deep expertise in crypto restructuring",
      "Flexible investment structures",
      "Focus on value recovery and turnaround situations"
    ],
    "investmentCriteria": [
      "Distressed but fundamentally sound projects",
      "Clear path to recovery or restructuring",
      "Attractive entry valuations",
      "Experienced management teams willing to restructure"
    ],
    "riskFactors": [
      "High risk of total loss",
      "Uncertain recovery timelines",
      "Legal and regulatory complications",
      "Market sentiment challenges"
    ],
    "expectedReturns": "Target IRR of 50-150% over 2-year period, with high risk/high reward profile",
    "portfolioComposition": "60% distressed tokens, 30% restructuring opportunities, 10% liquidation arbitrage"
  },
  "quickActions": [
    "Tell me about the team?",
    "What returns can I expect?",
    "What types of opportunities do you target?",
    "What are the fund terms?"
  ]
}
```

9. FUND DATA RETRIEVAL AUGMENTED GENERATION (RAG)

This use case focuses on venture capital funds, where the volume of documents per fund is relatively small (5–20 core files), including pitch decks, LP decks, side letters, and fund documentation. This constrained dataset simplifies retrieval design and makes it easy for an administrator to update files with ease. Different storage and context strategies can be applied depending on operational requirements, collaboration needs and the sensitivity of the data. Multiple approaches were evaluated:

1. BlueNexus - Account-Scoped Fund Context

BlueNexus's architecture is a comprehensive cloud platform combining secure data storage, search, memory, and execution environment for AI. Its distinguishing technical aspect is the deep integration of privacy tech (TEE, MPC) to ensure that personalization does not come at the cost of privacy or compliance. One can think of it as a secure data vault + app backend + AI memory all in one, which developers can integrate via SDK or API.

How it works:

- Each fund's documents are ingested into BlueNexus's secure vault.
- Documents are processed into vectorized embeddings, which are encrypted and stored in a per-account isolated environment.
- When a user queries the AI, BlueNexus retrieves only the minimal set of embeddings necessary to answer the question, keeping the underlying documents encrypted.
- The system leverages TEE and MPC to allow AI inference on sensitive fund data without ever exposing raw content.
- Access controls and audit logs are enforced at the vault level, ensuring compliance with LP confidentiality agreements and regulatory requirements.

This would be best suited for data that require strict access controls and highly auditable AI memory, where each fund or account must remain fully isolated. It makes it simple for administrators who want real-time updates to fund data while guaranteeing that AI queries never leak sensitive information.

2. XTrace - Encrypted Cross Agent Memory Layer

XTrace can be positioned as a portable, privacy-preserving memory layer that sits above models and applications, enabling persistent context across multiple AI agents and tools without re-exposing raw data.

How it works:

- User or team context (documents, conversations, preferences) is stored inside an encrypted vector memory hub.
- Memories are indexed for semantic search while remaining encrypted, using homomorphic techniques that allow retrieval without revealing plaintext.

- AI agents retrieve only the minimum relevant context at inference time via API, SDK, MCP, or browser extension - which is coming soon.
- Each user or organization operates within an isolated memory vault, ensuring strict separation of data.

XTrace is best suited for long-lived conversational and user-specific context. It is less optimized as a replacement for large, document-heavy RAG systems and instead complements them as a persistent memory layer.

3. Memvid - Portable Local Memory Layer

Memvid can be positioned as a local-first, single-file memory layer that enables fast, persistent knowledge retrieval for AI applications without requiring external databases or cloud infrastructure. It is designed to act as an embedded memory engine rather than a hosted service, prioritising portability, speed and data ownership.

How it works:

- Documents and text are ingested and converted into a single self-contained .mv2 file.
- Content is chunked, embedded and indexed using a hybrid approach (semantic vectors + lexical BM25).
- All data, embeddings, indexes and metadata are stored together in an append-only memory file inspired by video encoding.
- At query time, Memvid performs ultra-fast local retrieval to identify the most relevant memory frames.
- Retrieved text is decoded and passed to the model as grounded context for inference.

Each .mv2 file acts as an isolated memory capsule, which can represent a project, fund, agent, or knowledge domain. Files can be copied, versioned, shared, or deployed across environments with deterministic behaviour. nMemvid is best suited for document-centric RAG, long-term knowledge storage, and offline or on-prem AI systems. It is not a cross-user or multi-tenant memory service by default.

4. Space and Time - Verifiable Data Storage & Analytics (PostgreSQL)

Space and Time provides a verifiable data layer built on SQL (PostgreSQL-compatible) with cryptographic proofs over query execution. This approach is designed for funds that require structured, auditable, and queryable data, moving beyond traditional document-level Q&A into analytics with integrity guarantees.

How it works:

- Fund documents are parsed and normalized into structured tables.
- Data is stored in Space and Time's PostgreSQL environment.
- Queries executed are returned with results which have cryptographic verification.
- Outputs can be independently validated for integrity and correctness.

This approach is best suited for funds that want to move beyond document-level Q&A into structured analytics. Space and Time enables verifiable analytics at scale, making it appropriate for multi-fund platforms, institutional reporting, or environments where query correctness and verifiability are critical. Use cases such as portfolio analytics, performance dashboards, or investor reporting where accuracy and auditability are non-negotiable.

Implementation Applied

For this proof of concept, a custom backend RAG pipeline was implemented to demonstrate the core mechanics:

- Load and parse all PDFs from a fund's folder.
- Extract text using pdf-parse.
- Split text into 600-character chunks with overlap.
- Score chunks via keyword matching against the user query.
- Return the top 5 most relevant chunks into the confidential model's context.

The backend formats the prompt like this:

IMPORTANT: Answer using ONLY this information:

[Document Excerpt 1] ...

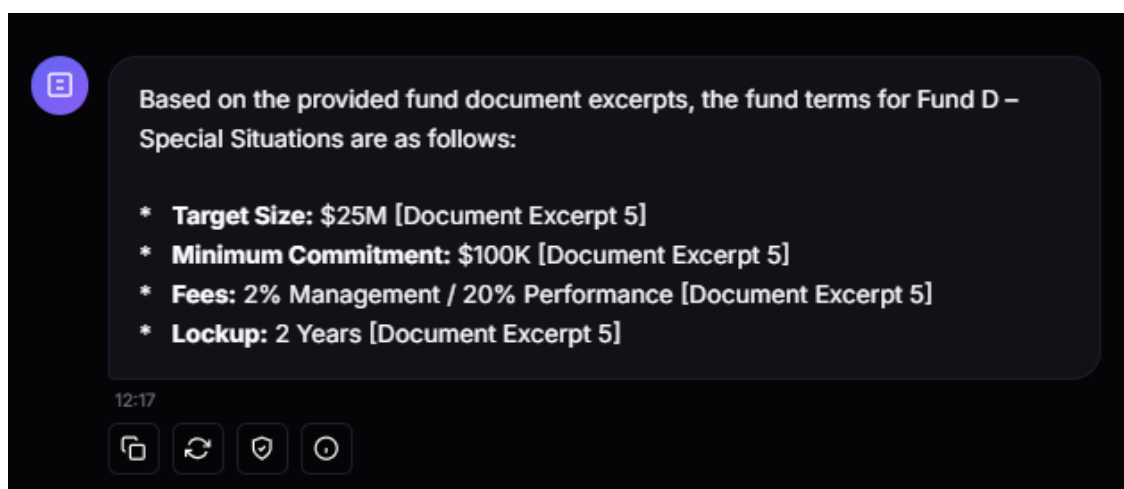
[Document Excerpt 2] ...

USER QUESTION: What is the fund size?

If the answer does not exist in the knowledge base, the model replies:

"I cannot answer this question as the information is not available in the provided knowledge base."

The prompt in the example below is: "What are the fund terms?"



Similar to the chat experience, LPs can independently verify both the message and the underlying model through the privacy verification centre. This enables recipients to confirm that responses were generated within a confidential execution environment, without exposing any underlying data. For a full product demonstration, please refer to the demo video [here](#).

10. DECENTRALIZED COMPONENTS RECAP

We've walked through the product flow and demo; now let's recap the core decentralized components. Each layer enforces trust, privacy, and integrity by design. Together, they form a cohesive, enterprise-grade AI stack: cloud-like in usability, but built on a foundation of verifiable, user-owned control. Every interaction, every query, every dataset is protected, auditable, and portable without compromise.

1. Confidential Compute

Purpose: Secure application and file hosting. Confidential compute provides a hardware-attested environment for code execution and document management, giving users full control over every file and line of code. By removing dependence on third-party cloud operators, it ensures that sensitive workflows are verifiably isolated from external access.

Implementation: Backend code and fund PDFs hosted on Phala CVM.

2. Confidential GPU TEE Models

Purpose: GPU-based TEEs allow AI models to process documents and context with cryptographic attestation. This enables multiple models to operate under a unified, auditable trust framework. Users can independently verify that every inference executed correctly, without exposing raw data or memory.

Implementation: GPU TEE Models provided via API through Redpill (Phala) paired with Privacy Verification Centre for independent attestation.

3. Sovereign Context

Purpose: Persistent, user-owned, portable AI memory. Sovereign context keeps chat history, uploaded documents, and other knowledge encrypted and under the user's control. Layers like XTrace, Memvid, BlueNexus, and Irys allow full ownership: context can be shared, revoked, or migrated without any operator ever gaining access.

Implementation: Login via BlueNexus (oAuth), export to XTrace Memory Vault, or upload to decentralized storage (Arweave via Irys).

4. Data Lakehouse / Warehouse

Purpose: Cryptographically verifiable structured storage with zero-knowledge guarantees. Documents are parsed, normalized, and stored in a SQL environment (Space and Time). Queries return results with ZK proofs, verifying correctness without exposing underlying data. This supports analytics, RAG, and structured reporting at enterprise scale.

Implementation: Create purpose-built "GPTs" using SxT Dreamspace, enabling shareable data pipelines with ZK-proven SQL, or create a Mv2 file to share context via portable file (Memvid).

11. KEY TAKEAWAYS

1. Tools are outpacing trust.

AI capabilities are advancing on a daily basis but the trust frameworks lag far behind. Early adopters may tolerate this imbalance but mainstream and institutional adoption - where sensitive data, third parties and accountability will expose the limits of "trust-me" verification. As systems become more autonomous and agentic, interacting with unknown counterparties and making consequential decisions, infrastructure must be able to verify identity, permissions, execution and outcomes without relying on a single controlling owner.

2. Centralized AI embeds unavoidable custodial risk.

Enterprise AI today depends on opaque infrastructure, contractual assurances and policy-based promises of non-retention. Every prompt, document and contextual signal processed by a centralized system remains an unprovable act of trust. Without cryptographic guarantees of non-access and correct execution, this model cannot satisfy the privacy, compliance or audit requirements of regulated and mission-critical environments.

3. The bottleneck is both the application and trust layer, not the model.

Open-source LLMs improve transparency but still fall short on usability, performance, and collaboration. Current centralized platforms dominate not because of superior models, but because they provide polished interfaces, persistent memory, expanded integrations and agent-like workflows. Any viable alternative must match these standards while fundamentally relocating trust away from operators and into verifiable hardware and cryptographic enforcement.

4. Verifiable privacy and data sovereignty are prerequisites, not trade-offs.

Decentralized AI enables enterprises to retain full ownership of data and model context without sacrificing cloud-like performance or collaboration. By combining confidential compute, GPU TEEs, sovereign memory and verifiable RAG pipelines, privacy, usability and auditability become complementary rather than conflicting goals. The real inflection point will come when users interact with AI seamlessly, unaware that every operation is executing inside a fully confidential, cryptographically enforced environment.

5. The shift is inevitable and operational, not ideological.

The move toward confidential AI is driven by operational necessity, not philosophy. Organizations handling high-value or sensitive data require systems where intelligence is delivered online, context is portable and user-owned and every interaction is provably correct and private. This architecture makes advanced AI not only possible, but practical and compliant for institutional-scale deployment.