# Generative Document Representation: A Research Landscape Analysis for a Semantic Agent Protocol

### **Executive Summary**

This report provides a comprehensive analysis of the academic and industry research landscape surrounding the concept of a Generative Document Representation (GDR) for a Semantic Agent Protocol. The GDR proposes a novel communication paradigm where autonomous AI agents exchange a compact Large Language Model (LLM) prompt as a semantic representation of a complex document or task. The receiving agent then uses this GDR with a specified, deterministic LLM to regenerate the identical, full-body document or execute the intended task. This analysis is structured around the five key technological pillars underpinning the GDR concept: Semantic Compression, AI Agent Communication Protocols, Decentralised AI (DeAI), Verifiable Computation, and the use of an LLM Prompt as a Formal Specification.

The analysis reveals that while the foundational technologies required for the GDR protocol are rapidly maturing, the concept itself represents a significant and novel leap beyond the current state of the art. In the domain of **semantic compression**, a notable maturity gap exists between the well-established field of Generative Visual Compression and the nascent area of generative text representation. Current research in text compression focuses primarily on shortening prompts to fit within context windows, whereas the GDR aims to use a prompt as a complete generative source for a document. Information theory provides a theoretical basis for this approach, suggesting that the semantic essence of a complex output can be encoded in a compact set of high-information tokens, though this implies a high degree of sensitivity in the GDR's construction.

The review of **AI agent communication protocols** indicates a clear industry trend toward modular, stacked protocols (MCP, A2A, ACP) that handle distinct layers of interaction such as tool use, discovery, and orchestration. The GDR protocol is best positioned not as a replacement for these standards, but as a novel, high-level *payload protocol* that defines the semantic content being exchanged. This positioning is critical for integration and adoption

within the emerging agentic ecosystem.

The **Decentralised AI (DeAI) ecosystem** provides a viable "full stack" for a trustless GDR marketplace. Projects focused on decentralized compute, agent frameworks, and on-chain governance offer the necessary components. The primary challenge lies in the integration and interoperability of these disparate systems. The concept of the Autonomous Economic Agent (AEA) creates a powerful economic incentive for a GDR marketplace, establishing a feedback loop where GDRs function as tradable, profit-generating digital assets.

**Verifiable computation** is essential for establishing trust in this decentralized marketplace. While Zero-Knowledge Machine Learning (zkML) offers the strongest cryptographic guarantees, its current computational overhead for large-scale LLM inference is prohibitive for most applications. This necessitates a "tiered verification" system, where the level of cryptographic proof is tied to the economic value of the task, blending zkML with more practical methods like Trusted Execution Environments (TEEs) and optimistic verification.

Finally, a critical analysis of the **LLM prompt as a formal specification** uncovers two fundamental challenges. First, the "determinism fallacy": setting an LLM's temperature to 0 does not guarantee bit-for-bit identical outputs due to the inherent nature of parallelized floating-point computations and modern Mixture-of-Experts (MoE) architectures. Second, the "prompt as contract" metaphor is technically and legally untenable due to issues of natural language ambiguity and security vulnerabilities.

Based on these findings, this report concludes with strategic recommendations. The GDR protocol must be architected to be resilient to non-determinism, pivoting from verifying "identical regeneration" to "semantic equivalence." The prohibitive cost of verification must be managed through a tiered, value-based system. Most critically, the conceptual framing should shift from "prompt as contract" to "prompt as blueprint," where a robust on-chain smart contract serves as the legally and technically binding agreement governing the execution of the generative instructions contained within the GDR.

# I. Introduction: The Emerging Paradigm of Semantic, Agent-Driven Communication

The contemporary digital ecosystem is undergoing a fundamental transformation, shifting from rigid, human-driven interfaces to dynamic, autonomous systems powered by artificial intelligence. This transition necessitates a parallel evolution in communication protocols. For decades, machine-to-machine interaction has been dominated by syntactic, API-driven models, where data is exchanged in highly structured and inflexible formats like JSON or XML.

While effective for predictable, programmatic tasks, these models are ill-suited for the fluid, context-aware, and intent-driven interactions that characterize the emerging world of AI agents. The growing complexity of tasks being delegated to autonomous systems reveals the brittleness of these legacy protocols, requiring bespoke connectors and constant maintenance to integrate fragmented systems.

In this context, the Generative Document Representation (GDR) protocol is proposed as a forward-looking alternative. The core proposition is to leverage the expressive power of Large Language Models (LLMs) to create a new communication primitive based on semantic intent rather than syntactic structure. The GDR is a highly compact, self-contained LLM prompt that functions as a "semantic hash" or generative source code for a complex document or task. Instead of exchanging a multi-megabyte data file, autonomous agents would exchange a succinct GDR, measured in kilobytes. The receiving agent, using a contractually specified and deterministic LLM (e.g., a specific model hash at temperature=0), can then perfectly regenerate the full-body document or execute the intended task. This paradigm moves agent communication from a declarative model of exchanging data to a generative model of exchanging instructions, promising orders-of-magnitude improvements in efficiency and flexibility.

This report provides a comprehensive research landscape analysis to evaluate the feasibility, novelty, and challenges of the GDR concept. The objective is to map the existing body of academic and industry research to the five foundational pillars upon which the GDR protocol is built. These pillars are:

- 1. **Generative and Semantic Compression:** The theoretical and practical underpinnings of using generative models to create compact representations of complex data.
- 2. **Al Agent Communication Protocols:** The evolving standards and architectures that govern how autonomous agents interact, discover, and collaborate.
- 3. **The Decentralised AI Ecosystem:** The Web3 infrastructure, including marketplaces for compute and AI services, on-chain agents, and DAOs, that can support a trustless, open economic system for GDR-based tasks.
- 4. **Verifiable Computation and zkML:** The cryptographic techniques that can provide trust, privacy, and auditability for agent actions in a decentralized environment.
- 5. **The LLM Prompt as a Formal Specification:** A critical examination of the reliability, determinism, and security of using natural language prompts as unambiguous, contract-like instructions.

By systematically reviewing the state-of-the-art in each of these domains, this report aims to serve as a foundational document for the GDR project's strategic planning and R&D roadmap. It will identify areas where the GDR concept aligns with established research, pinpoint its most novel contributions, and, most importantly, highlight the critical technical hurdles and conceptual fallacies that must be addressed for its successful implementation.

# II. Pillar 1: The State of Generative and Semantic Compression

The foundational premise of the Generative Document Representation (GDR) is its ability to achieve extreme compression by storing generative instructions rather than the final output. This section critically examines this concept by first establishing its theoretical basis in information theory, then drawing powerful analogies from the mature field of generative visual compression, and finally surveying the nascent but directly relevant research in semantic text compression.

### 2.1 Foundations in Information Theory and LLMs

The connection between information theory and language modeling dates back to Claude Shannon's foundational work, which measured how well simple n-gram models could predict or compress natural language text.<sup>3</sup> Modern Large Language Models (LLMs) can be viewed as the ultimate evolution of this principle; they are fundamentally powerful compression and prediction engines, trained to model the statistical distribution of human language.<sup>3</sup> Their ability to generate coherent text is a direct consequence of their ability to predict the next token in a sequence, a process that inherently involves compressing vast amounts of training data into a compact set of model parameters.

Recent research has begun to apply information-theoretic principles more directly to analyze the internal mechanisms of LLMs. One key concept is the **information bottleneck (IB) principle**, which characterizes the trade-off between compressing input information and preserving the predictive power necessary for a given task. In the context of multimodal models, for instance, researchers have used the IB perspective to understand "visual forgetting," where excessive compression of visual representations during text-driven instruction tuning leads to a loss of crucial knowledge. This principle is directly relevant to GDR's core challenge: achieving maximal compression of a document's semantic essence into a prompt without losing the necessary information for its faithful regeneration.

Further supporting the feasibility of GDR is the discovery of "mutual information (MI) peaks" in the reasoning trajectories of LLMs.<sup>6</sup> Studies show that during a complex reasoning process, the mutual information between the model's internal state and the correct final answer does not increase smoothly. Instead, it exhibits sudden, significant spikes at specific, sparse

generative steps. These MI peaks often correspond to "thinking tokens" such as "Hmm," "Wait," or "Therefore," which signal reflection or logical transition. Suppressing the generation of these tokens has been shown to degrade reasoning performance, while other tokens have minimal impact. This phenomenon suggests that the core semantic and logical structure of a complex output is not evenly distributed but is concentrated in a few critical nodes. This provides strong theoretical support for the GDR hypothesis that a compact set of instructions, if correctly formulated to target these high-information nodes, can effectively encapsulate the generative instructions for a much larger and more complex output.

### 2.2 Generative Compression: Lessons from the Visual Domain

While generative compression for text is an emerging field, the domain of visual data (images and video) offers a rich and mature body of research that serves as a powerful analogue for the GDR concept. Generative Visual Compression (GVC) moves beyond traditional codecs like JPEG or H.266/VVC, which focus on eliminating statistical redundancy in the pixel space (e.g., by discarding high-frequency details). Instead, GVC leverages deep generative models—such as Variational Auto-Encoders (VAEs), Generative Adversarial Networks (GANs), and Diffusion Models (DMs)—to learn the underlying distribution of visual data. The core principle is to encode an image into a compact latent representation and then use a powerful generative decoder to synthesize a visually pleasing reconstruction. This "analysis-by-synthesis" approach is conceptually identical to GDR's "prompt-to-document" methodology.

The advantages of GVC are directly translatable to the goals of GDR:

- Compact Feature Representation: Generative models excel at learning compact, semantically rich feature distributions, enabling far higher compression ratios than traditional methods, especially at ultra-low bitrates.<sup>7</sup>
- **High-Fidelity Reconstruction:** By learning the "prior" of what natural images look like, generative decoders can produce photo-realistic reconstructions that are perceptually convincing, even if they are not pixel-perfect copies of the original. This aligns with GDR's goal of semantic, rather than merely syntactic, fidelity.

Recent advancements in GVC offer specific technical parallels. For example, the **Control-GIC** framework uses a VQGAN to represent an image as a variable-length sequence of codes, enabling fine-grained bitrate adaptation—a feature that could be analogous to a GDR system with varying levels of detail.<sup>15</sup> The **Generative Latent Coding (GLC)** architecture performs compression in the latent space of a VQ-VAE, which is sparser and more semantically aligned with human perception, achieving high-realism compression at less than 0.04 bits per pixel (bpp).<sup>11</sup>

Most relevant to GDR is the rise of **multimodal semantic compression**. Recent work introduces frameworks like Multimodal Image Semantic Compression (MISC), which stores images using key descriptive elements and context—often derived from a textual prompt—instead of pixel-level details. This approach uses text-to-image models like Stable Diffusion to generate high-quality visuals from simple textual prompts, demonstrating a direct link between a compact semantic input (text) and a complex, high-fidelity output (image). This is a direct and powerful proof-of-concept for the GDR paradigm in a different modality. The visual domain has clearly established that compressing the semantic essence and using a generative model for reconstruction is a viable and highly effective strategy.

# 2.3 Semantic Compression for Text: From Context Reduction to Document Representation

The application of generative and semantic compression to text is a more nascent field, but one that is rapidly gaining traction, driven largely by the constraints of LLM context windows. Traditional semantic compression in natural language processing involved compacting a lexicon by replacing less frequent terms with their more general hypernyms, a process aimed at reducing dimensionality for information retrieval tasks.<sup>17</sup> Modern approaches, however, leverage LLMs themselves as the compression engine.

Current research in LLM-based text compression can be broadly divided into two categories <sup>18</sup>:

- 1. **Selective Compression:** These methods aim to identify and retain the most important tokens or sentences from an original prompt while discarding the rest. The goal is to shorten the input context while preserving key information. However, this approach risks information loss, as discarded tokens are never compensated for.<sup>19</sup>
- 2. **Generative Compression:** These methods utilize a language model to take an original prompt as input and generate a shorter, more concise version. This approach, which involves rewriting or summarization, can better preserve contextual integrity and global structure but may suffer from hallucination.<sup>18</sup>

A prominent example of the generative approach is the **SCOPE** framework, which proposes a "chunking-and-summarization" mechanism.<sup>19</sup> SCOPE splits a long prompt into semantically coherent chunks, evaluates each chunk's relevance, and then uses a small summarization model to rewrite the chunks, prioritizing the compression of less relevant or longer chunks. A key innovation is its dynamic compression ratio control, which calculates a specific compression target for each chunk based on its length () and relevance (), ensuring that more important information is preserved with higher fidelity.<sup>19</sup> This demonstrates a sophisticated method for optimizing the trade-off between compression and information preservation, a

technique that could be adapted for creating optimized GDRs.

The most foundational work in this area is the paper "Semantic Compression With Large Language Models" by Gilbert et al..<sup>21</sup> This research directly explores the viability of using LLMs like GPT-4 to perform their own "approximate compression" on prompts. The methodology involves instructing an LLM to compress a piece of text and then, in a separate session, instructing it to decompress the output. Their initial results are highly promising, indicating that GPT-4 can effectively compress and reconstruct text while preserving its semantic essence, achieving a compression ratio of approximately 5x over the original token count.<sup>22</sup> Crucially, they introduce two novel metrics for evaluating the quality of this process:

- Exact Reconstructive Effectiveness (ERE): Measures the degree of perfect, character-for-character reconstruction.
- Semantic Reconstruction Effectiveness (SRE): Measures the preservation of intent and meaning, typically using cosine similarity between the embeddings of the original and reconstructed text.

These metrics provide a direct and practical framework for evaluating the fidelity of GDR regeneration. Further work by Fei et al. draws inspiration from source coding in information theory, using a pre-trained model to reduce the semantic redundancy of long inputs before passing them to an LLM, enabling generalization to texts that are 6-8 times longer without fine-tuning.<sup>24</sup>

While this body of research validates the core principle of using an LLM for compression and regeneration, it also highlights a critical distinction. The primary motivation for nearly all current work is to shorten long contexts to fit within an LLM's input window. The GDR concept is more ambitious. It proposes using the prompt not merely as a compressed version of an existing document, but as the canonical, generative representation of a document that may not have existed before. This moves beyond mere compression to a paradigm of generative representation. The existing research provides a strong foundation and valuable tools (like SCOPE's optimization and SRE metrics), but GDR's application of these ideas represents a novel and significant extension of the current state of the art.

# III. Pillar 2: The Evolving Landscape of Al Agent Communication

For the Generative Document Representation (GDR) to function as an effective medium of exchange, it must operate within a coherent and widely adopted communication framework. The history of AI agent communication reveals a clear trajectory from monolithic, semantically rigid standards to a more flexible, modular ecosystem. This section analyzes this evolution,

details the current state-of-the-art in agent communication protocols, and positions the GDR concept within this modern landscape.

#### 3.1 From Monolithic Standards to Modular Protocol Stacks

The foundational work in agent communication was pioneered by standards like the **Knowledge Query and Manipulation Language (KQML)** and, more formally, the **FIPA Agent Communication Language (FIPA-ACL)**.<sup>25</sup> These protocols were revolutionary in that they were based on speech act theory, treating inter-agent messages not as simple data packets but as intentional communicative acts, or "performatives".<sup>27</sup>

The FIPA-ACL specification defines a structured message format with a set of parameters to describe the context and content of the communication.<sup>29</sup> The only mandatory parameter is the performative, which denotes the type of communicative act (e.g., inform, request, query-if, propose). Other key parameters include:

- Participants: sender, receiver, reply-to.
- Content: content, language, encoding, ontology.
- Conversation Control: protocol, conversation-id, in-reply-to.

This structure allowed for complex, patterned conversations, such as negotiations and auctions, to be standardized.<sup>27</sup> A critical aspect of these early protocols was their reliance on a shared ontology—a formal, explicit specification of concepts and their relationships—to ensure that communicating agents had a common, unambiguous understanding of the terms used in the message content.<sup>26</sup> While powerful, this requirement for pre-defined, shared ontologies proved to be a bottleneck, limiting flexibility and scalability in open, heterogeneous environments.

The recent explosion of LLM-based agentic AI has led to a paradigm shift. The focus has moved away from a single, all-encompassing standard like FIPA-ACL towards a more pragmatic and modular "protocol stack". In this model, different protocols are layered to handle distinct aspects of agent interaction, mirroring the layered architecture of the internet (e.g., TCP/IP for transport, HTTP for application). This modularity allows for greater flexibility and specialization, enabling the rapid development and integration of diverse agent systems.

### 3.2 Comparative Analysis of Modern Protocols: MCP, A2A, and ACP

The contemporary agent communication landscape is dominated by three emerging, complementary protocols, each addressing a different layer of the interaction stack.

- 1. **Model Context Protocol (MCP):** Developed by Anthropic and now an open-source standard, MCP is designed to standardize how AI agents connect to and interact with external systems, such as APIs, databases, local files, and other tools.<sup>2</sup> It functions as a universal "toolbelt" for agents. MCP operates on a client-server architecture where an MCP Host can manage multiple MCP Clients, each maintaining a one-to-one connection with an MCP Server that exposes a tool or data source.<sup>33</sup> Communication occurs via a simple set of message types (Request, Result, Error, Notification) formatted as JSON-RPC 2.0, transmitted over transports like stdio or Server-Sent Events (SSE).<sup>1</sup>
- 2. **Agent-to-Agent (A2A) Protocol:** Spearheaded by Google and managed by the Linux Foundation, the A2A protocol focuses on inter-agent discovery, negotiation, and collaboration.<sup>2</sup> It acts as the "social network" for agents, enabling them to find each other and coordinate tasks. A key innovation is the concept of "**Agent Cards,"** which are self-descriptions that agents publish to advertise their capabilities, supported protocols, and accepted request formats.<sup>32</sup> This allows for dynamic discovery and matchmaking in an open ecosystem. A2A follows a client-server model for communication over HTTPS, using JSON messages for data exchange.<sup>1</sup>
- 3. **Agent Communication Protocol (ACP):** An initiative led by IBM and also part of the Linux Foundation, ACP is designed for orchestrating complex workflows, delegating tasks, and maintaining state across multiple agents. It serves as the "project manager" of a multi-agent system. ACP is built on a REST-first, HTTP-native architecture, making it easy to integrate into existing enterprise environments using standard tools. It is designed to be async-first, which is ideal for long-running tasks, and supports offline discovery, where agents can be found even when inactive.

These protocols are not competitors but complementary components of a comprehensive agent communication stack. An agent might use A2A to discover a specialized service agent, use ACP to delegate a multi-step task to it, and that service agent might in turn use MCP to access the specific tools and data needed to complete the task.<sup>2</sup>

Table 1: Comparative Analysis of Modern AI Agent Communication Protocols

Dimension	FIPA-ACL (Baseline)	Model Context Protocol (MCP)	Agent-to-Age nt (A2A) Protocol	Agent Communicati on Protocol (ACP)
Primary Focus	Semantic Intent & Structured	Tool/Resource Access	Agent Discovery &	Workflow Orchestration & State

	Dialogue		Collaboration	Management
Core Abstraction	"Speech Act"	"Universal Toolbelt"	"Social Network"	"Project Manager"
Architecture	Message-base d (transport agnostic)	Client-Server	Client-Server	Client-Server (RESTful)
Communicati on Style	Structured Message Parameters	JSON-RPC 2.0 over stdio/SSE	JSON over HTTPS	RESTful API over HTTP
Key Features	Performatives, Shared Ontology	One-to-one client-server relationship, Standardized tool use	"Agent Cards" for discovery, Capability negotiation	Async-first, Offline discovery, Multimodal messages
Role in GDR Ecosystem	Provides semantic framework for the act of sending a GDR.	Provides access to the LLM specified in the GDR for execution.	Enables discovery of agents capable of creating or executing GDRs.	Transports GDR packets as payloads within a managed workflow.

# 3.3 Semantic Negotiation and Ontology Mediation for True Interoperability

The modern protocol stack effectively addresses the mechanics of connection, discovery, and orchestration. However, it largely sidesteps the original challenge that FIPA-ACL sought to solve: ensuring a shared semantic understanding of the *content* being communicated. Even if two agents can successfully exchange messages via ACP, they may fail to collaborate if their internal representations of concepts (their ontologies) differ. For example, one agent's "customer\_id" might be another's "client\_ref".

Research in this area provides two main approaches to bridge this semantic gap:

- Ontology Mediation Methods: These are techniques for creating mappings between
  different, static ontologies.<sup>37</sup> This includes ontology matching (finding correspondences
  like "price" ≈ "cost"), ontology alignment (defining formal rules like "Car is-a Vehicle"),
  and ontology merging (combining multiple ontologies into a single, unified one). These
  methods are effective when agent ontologies are known in advance but are less suited
  for dynamic, open environments.
- Ontology Negotiation Protocols (ONP): This more dynamic approach allows agents to
  discover and resolve semantic conflicts at runtime.<sup>37</sup> Based on the influential work of
  Bailin and Truszkowski, an ONP enables an agent that encounters an unfamiliar term to
  engage in a structured dialogue of interpretation, clarification, and explanation with its
  counterpart until a common understanding is reached or the conflict is deemed
  irreconcilable.<sup>37</sup> This is crucial for enabling effective collaboration between agents that
  were not pre-designed to work together.<sup>38</sup>

### 3.4 Implications for a GDR-Based Protocol

The analysis of the modern agent communication landscape provides a clear strategic path for the GDR protocol. It should not be positioned as a new transport or orchestration layer to compete with ACP or A2A. Instead, GDR is a novel **payload protocol** or **content type**. It defines a standard for encoding complex, generative tasks into a compact, semantic format.

In this model, a GDR packet—containing the prompt, specified LLM, reward, deadline, and verification requirements—becomes the content of a message transmitted via an established protocol. For example:

- 1. A *client agent* uses the A2A protocol to discover a network of *worker agents* that advertise their capability to execute GDRs for a certain class of LLMs.
- 2. The client agent initiates a task by sending a request message to a chosen worker agent using the ACP protocol.
- 3. The *payload* of this ACP message is a GDR packet.
- 4. The worker agent receives the GDR, uses MCP to connect to the specified LLM resource (e.g., an API or a decentralized compute node), and executes the generative task.
- 5. Upon completion, the worker agent generates a proof of execution and returns the result (or a pointer to it) to the client agent.

This approach allows the GDR protocol to leverage the rapidly growing ecosystem of agent communication infrastructure, focusing its innovation on the unique challenge of representing semantic intent, rather than reinventing the mechanics of agent interaction. It also presents a unique hybrid approach to semantic agreement. While the natural language prompt relies on

the "shallow" contextual understanding of the LLM, the protocol demands a "deep," cryptographically precise agreement on the execution environment (the exact model hash, temperature, etc.). This creates a new interoperability requirement not central to existing protocols: the ability for agents to negotiate and verify the precise computational context of a task.

# IV. Pillar 3: The Decentralized AI Ecosystem as a Foundation for GDR

The vision for the GDR protocol extends beyond simple agent-to-agent communication to encompass a trustless, open marketplace for generative tasks. This requires a robust, decentralized infrastructure that can support computation, data exchange, agent interaction, and governance without relying on a central authority. This section surveys the Decentralised AI (DeAI) and Web3 landscape, demonstrating that the components necessary for such a marketplace are not merely theoretical but are actively being developed and deployed.

### 4.1 The DeAl Marketplace: Compute, Data, and Model Economies

The Decentralised AI (DeAI) movement has emerged as a direct response to the increasing centralization of AI development within a handful of large technology corporations.<sup>39</sup> By leveraging blockchain and other distributed technologies, DeAI aims to democratize access to the core resources of AI: data, computational power, and models. This creates the ideal foundation for a GDR marketplace. The DeAI landscape can be categorized into several key functional layers:

- Decentralised Compute Marketplaces: These platforms create a peer-to-peer market for computational resources, particularly GPUs, which are essential for running LLMs. Projects like io.net <sup>41</sup>, Akash Network <sup>42</sup>, and Golem <sup>42</sup> allow individuals and data centers with underutilized hardware to rent out their capacity to those who need it, often at a fraction of the cost of centralized cloud providers. <sup>41</sup> These networks provide the physical execution layer where worker agents can run the LLM inference specified in a GDR.
- Decentralised AI Service and Model Marketplaces: These platforms facilitate the
  discovery, sharing, and monetization of AI models and services. SingularityNET aims to
  create a global marketplace where anyone can buy and sell AI algorithms and services.<sup>42</sup>
  Bittensor takes a different approach, creating an incentivized, peer-to-peer network
  where AI models compete and collaborate, with participants being rewarded in

- cryptocurrency for contributing valuable intelligence.<sup>42</sup> Such platforms could host a marketplace where GDRs themselves are traded, or where the specialized models needed to execute them are made available.
- **Decentralised Data Marketplaces:** Projects like **Ocean Protocol** enable secure and privacy-preserving data sharing and monetization. <sup>42</sup> By tokenizing datasets, Ocean creates a marketplace where data owners can sell access to their data without losing control over it. This is relevant for GDRs that may require access to specific, proprietary datasets for their execution.

While this ecosystem is vibrant, it faces significant challenges, including performance bottlenecks due to blockchain latency, scalability constraints, and the coordination overhead of managing distributed resources. <sup>46</sup> Nevertheless, the active development across these layers provides a clear path for sourcing the off-chain resources required by the GDR protocol in a decentralized manner.

Table 2: Overview of the Decentralized AI (DeAI) Marketplace Landscape

Project Name	Primary Function	Core Concept	Relevance to GDR Ecosystem
io.net / Akash Network	Decentralized GPU Compute	Creates a peer-to-peer marketplace for renting underutilized GPU power for AI/ML workloads.	Provides the distributed, cost-effective physical compute layer for worker agents to execute LLM inference as specified by a GDR.
Fetch.ai	Autonomous Economic Agent Platform	Provides a framework for building and deploying Autonomous Economic Agents (AEAs) that can transact and operate on-chain.	Provides the "worker agents" that would discover, acquire, and execute GDR tasks to generate economic value.
SingularityNET	Decentralized AI Service	Allows developers to create, share,	Could serve as a marketplace where

	Marketplace	and monetize AI services, algorithms, and models in a global, open market.	GDRs are listed as "jobs" or where specialized models required by GDRs are sourced.
Bittensor	Incentivized Peer-to-Peer ML Network	Creates a competitive market for intelligence where models are rewarded for providing value to the collective network.	Offers a mechanism for incentivizing the creation of high-quality models that could be used for GDR execution or verification.
Ocean Protocol	Decentralized Data Marketplace	Enables the tokenization of data assets, allowing for privacy-preserving data sharing and monetization.	Provides a secure way for GDR tasks to access specific, proprietary datasets without compromising data ownership.
Numerai	Crowdsourced AI Hedge Fund	Incentivizes data scientists to build predictive models on encrypted financial data, rewarding the best models with cryptocurrency.	Demonstrates a successful model of incentivizing a global network of specialists to perform complex computational tasks for a reward.

# 4.2 The On-Chain Al Agent: Architecture, Sovereignty, and Economic Models

Parallel to the development of DeAI infrastructure is the emergence of the "On-Chain AI

**Agent Economy"**. <sup>47</sup> This paradigm envisions AI agents not just as off-chain programs that interact with blockchains, but as first-class citizens within them: self-sovereign, autonomous entities capable of owning assets, executing transactions, and engaging in economic activity without direct human intervention. <sup>47</sup>

The architectural pillars required to realize this vision are 47:

- 1. **Sovereign Wallets:** Each agent requires its own blockchain wallet to hold funds, manage digital assets (like NFTs), and pay for transactions (gas fees), giving it true economic independence.
- 2. **Secure Runtime Environments:** A protected and verifiable environment where the agent's logic can execute without interference.
- 3. **Secure Access to Web and Chains:** Agents need the ability to access both off-chain data (via oracles) and on-chain data and services to make informed decisions.

This model transforms the role of the blockchain from a simple ledger for payments into a trust and identity layer for AI. Every action an agent takes can be recorded immutably, creating a verifiable and auditable history of its behavior. <sup>49</sup> Projects like **Fetch.ai** are specifically designed to support these **Autonomous Economic Agents (AEAs)**, providing the tools to build agents that can autonomously pursue economic goals. <sup>42</sup>

The GDR concept integrates seamlessly into this economic model. A GDR, representing a valuable task with a specified reward, can be minted as a Non-Fungible Token (NFT). This makes the task a unique, ownable, and tradable digital asset. An AEA could be programmed to scan a marketplace of GDR NFTs, analyze the potential profit of a task (reward minus estimated compute and gas costs), purchase the NFT, execute the task on a decentralized compute network, and collect the reward upon successful verification. This creates a dynamic, self-sustaining economy around GDRs, where the pursuit of profit by a network of competing AEAs drives the efficient execution of generative tasks.

### 4.3 Governance and Treasury Management via Al-Integrated DAOs

In a decentralized ecosystem, governance is paramount. **Decentralized Autonomous Organizations (DAOs)** provide a framework for community-led governance, where rules are encoded in smart contracts and decisions are made through token-holder voting. <sup>54</sup> DAOs are the natural choice for governing the GDR protocol itself, enabling the community to manage protocol upgrades, curate lists of certified and trusted LLM models, and oversee dispute resolution mechanisms.

Furthermore, there is a strong trend toward integrating AI directly into DAO operations to

enhance efficiency and intelligence.<sup>56</sup> Al agents are being developed to automate a wide range of DAO functions, including <sup>58</sup>:

- Proposal Analysis: Al can summarize complex governance proposals, analyze their potential impact on the treasury, and assess risks, helping token holders make more informed voting decisions.
- Treasury Management: All can monitor market conditions and execute pre-approved investment strategies to manage the DAO's treasury, optimizing yield and preserving capital.
- Automated Governance: For certain classes of decisions, AI agents can be delegated voting power to act on behalf of human members or even other DAOs, a concept known as "metagovernance."

This synergy between AI and DAOs creates a robust framework for managing a complex protocol like GDR. A DAO could govern the core protocol, while a swarm of specialized AI agents, under the DAO's purview, could manage the day-to-day operations of the marketplace, such as monitoring agent reputation, flagging malicious activity, and optimizing the allocation of tasks to compute providers. This combination of decentralized human oversight and autonomous AI execution provides a scalable and resilient model for the GDR ecosystem.

# V. Pillar 4: Trust and Privacy through Verifiable Computation

For a decentralized marketplace of GDR-based tasks to function, participants need a way to trust each other without relying on a central intermediary. A task creator needs assurance that a worker agent has correctly executed the specified GDR using the correct LLM, and the worker agent needs a way to prove its work to receive payment. Similarly, if the GDR or its output contains sensitive information, privacy must be maintained. Verifiable computation, particularly Zero-Knowledge Machine Learning (zkML), provides the cryptographic foundation for these guarantees.

### 5.1 Principles of Zero-Knowledge Machine Learning (zkML)

Zero-Knowledge Machine Learning (zkML) is a powerful cryptographic technique that merges Zero-Knowledge Proofs (ZKPs) with machine learning computations.<sup>59</sup> A ZKP allows a "prover"

to convince a "verifier" that a statement is true without revealing any information beyond the validity of the statement itself. In the context of zkML, the prover can generate a succinct cryptographic proof attesting to the correct execution of an ML model's inference.

The process generally involves three main steps <sup>59</sup>:

- 1. **Compilation:** The ML model's computational graph (e.g., the layers of a neural network) is converted into a mathematical representation, typically an **arithmetic circuit** or a system of constraints.
- 2. **Proof Generation:** The prover executes the model inference and, in parallel, uses the circuit and the input data to generate a compact ZK proof (such as a zk-SNARK). This proof cryptographically binds the inputs, the model, and the output, attesting that the computation was performed correctly.
- 3. **Verification:** The verifier, which can be a resource-constrained entity like a blockchain smart contract, can then check the validity of this proof very quickly and cheaply, without needing to re-run the entire computation or have access to the private inputs or the model's weights.

The primary role of zkML in a blockchain context is to enable the **trustless verification of computationally intensive off-chain work**. Since LLM inference is far too complex to be executed directly on a blockchain, it must happen off-chain. zkML provides the bridge, allowing an on-chain smart contract to verify with mathematical certainty that the off-chain computation was performed correctly, thus enabling trustless interactions between agents. <sup>62</sup>

#### 5.2 Architectures for Verifiable LLM Inference

While zkML offers the strongest, purely cryptographic guarantees, it is not the only method for achieving verifiable inference. The field is exploring several approaches, each with distinct trade-offs in terms of trust assumptions, performance, and cost.<sup>63</sup>

- 1. **Zero-Knowledge Proofs (zkML):** This software-based approach relies solely on mathematical and cryptographic assumptions. It offers the highest level of trustlessness and can provide strong privacy for both inputs and model weights. Its primary drawback is the extremely high computational overhead required to generate the proof.<sup>60</sup>
- 2. Trusted Execution Environments (TEEs): This hardware-based approach utilizes secure enclaves within a processor (like Intel SGX or AMD SEV) to run computations in an isolated, encrypted environment. A remote attestation process allows a user to verify that their code is running untampered inside a genuine TEE. TEEs offer very low performance overhead but require trusting the hardware manufacturer to have designed and fabricated the TEE securely.<sup>63</sup>
- 3. Optimistic Verification (OPML): This game-theoretic approach, inspired by optimistic

rollups in blockchain scaling, assumes computations are correct by default. The result is posted on-chain, and a "challenge period" opens, during which any independent validator can re-execute the computation and submit a "fraud proof" if they find a discrepancy. This method is computationally cheap in the happy path but introduces significant latency (the length of the challenge period) and is less effective for non-deterministic computations.<sup>63</sup>

4. **Hashing and Fingerprinting:** This probabilistic approach involves the prover generating hashes of intermediate computational states (e.g., the outputs of certain layers) during inference. A verifier can then randomly select a few steps, re-compute them, and check if the hashes match. This adds very low overhead but provides statistical, rather than absolute, guarantees of correctness.<sup>63</sup>

The existence of these diverse methods suggests that a one-size-fits-all approach to verification may not be optimal. The economic value and security requirements of a given GDR task should dictate the appropriate level of verification. For a low-value, public data generation task, a probabilistic check or even just agent reputation might suffice. For a high-value task involving private financial data, a full ZK proof or a TEE-based execution would be necessary. This points toward the need for a "tiered verification" system where the cost of trust is commensurate with the value at stake.

Table 3: Trade-offs in Verifiable Inference Methodologies

Methodolo gy	Trust Assumptio n	Computati onal Overhead (Prover)	Verificatio n Latency	Privacy Guarantee s	Suitability for GDR
zkML	Cryptograp hic/Mathem atical	Extremely High (100-10,00 Ox)	Near-instan t (on-chain)	High (can hide inputs/mod el)	"Gold standard" for high-value/ privacy tasks but cost-prohib itive for most.
TEEs	Hardware Manufactur er	Very Low	Near-instan t	High (hides from host OS/infra)	Excellent for medium-to -high value, latency-sen

					sitive tasks where hardware trust is acceptable.
Optimistic Verificatio n	Game Theoretic (≥1 honest verifier)	Low (re-comput ation on challenge)	High (challenge period, e.g., days)	None (data must be public for verification)	Unsuitable for real-time tasks; viable for asynchrono us, low-value tasks with public data.
Hashing/Fi ngerprintin g	Probabilisti c	Low (~5-20%)	Near-instan t	None	Good for low-trust audit trails and detecting gross negligence, but not for cryptograp hic security.

### 5.3 Analysis of Computational Overhead and Practicality for GDR

The primary obstacle to the widespread adoption of zkML for verifying LLM inference is its practicality, specifically the immense computational overhead. Research consistently highlights the "significant inefficiency and tremendous overhead" associated with generating ZK proofs for the complex, non-linear operations found in modern neural networks.<sup>64</sup>

General-purpose ZK-SNARK systems, while flexible, are not optimized for ML workloads. As a result, proving a single inference for a large model like Llama2-7B could take **"hundreds of hours"** and require massive amounts of memory, making the approach impractical for all but

the most trivial models.<sup>65</sup> The cost of generating the proof would vastly exceed the cost of the original inference itself, rendering any economic model based on it unviable.

In response, a new wave of research is focused on developing specialized protocols and frameworks to accelerate ZKML.

- zkGPT proposes optimizations like constraint fusion and circuit squeezing tailored for transformer architectures to reduce the overhead of proving non-linear layers like GeLU and attention.<sup>64</sup>
- **ZKTorch** is an end-to-end system that compiles ML models into compositions of optimized cryptographic primitives. By extending the Mira accumulation scheme to allow for parallel proof aggregation, ZKTorch claims to achieve up to a **6x speedup in proving time** and a **3x reduction in proof size** compared to previous specialized protocols.<sup>66</sup>

Despite these advances, the overhead remains substantial. This reinforces the conclusion that a full ZK proof is not a practical solution for every GDR task. A tiered verification system, allowing task creators to select a verification method appropriate for their budget and security needs, is an architectural necessity for the GDR marketplace to be economically sustainable. Furthermore, a robust GDR protocol must verify not only the correctness of the computation but also that the *correct model* was used. A malicious provider has a strong incentive to substitute a smaller, cheaper model.<sup>67</sup> Therefore, the verification process must incorporate a commitment to the model's identity, for example, by including a hash of the model weights in the public inputs of the ZK proof.

### 5.4 The Commercial and Open-Source Landscape for Verifiable Al

The critical need for trust in AI is driving the growth of a commercial and open-source ecosystem dedicated to verifiable inference.

- Commercial Platforms: Companies are emerging to productize these complex cryptographic techniques. EigenCloud has launched EigenAI, a verifiable LLM inference API that aims to provide unaltered prompts and responses, and EigenCompute, a secure off-chain execution environment using TEEs with a roadmap toward ZK proofs.<sup>68</sup>
   Hyperbolic is building a decentralized compute network that incorporates a "Proof of Sampling" protocol to cryptographically verify that model outputs are genuine.<sup>70</sup>
- Open-Source Projects: The research community is actively developing open-source solutions. The SVIP protocol is a novel, lightweight method for verifiable inference that does not rely on heavy cryptography. Instead, it requires the provider to return processed hidden representations from the LLM, which are used to train a proxy task that acts as a unique model identifier, allowing users to detect model substitution with very low overhead.<sup>67</sup> Parity Protocol is another open-source project building a decentralized

compute engine where tasks are executed across multiple nodes, and outputs are matched and verified before acceptance.<sup>71</sup>

This burgeoning landscape of both commercial APIs and open-source tools provides a rich set of options for implementing the verification layer of the GDR protocol. A strategic approach would involve building an abstraction layer that can integrate with multiple verification backends (zkML, TEE, optimistic, etc.), allowing the GDR marketplace to adapt as these technologies mature and their cost-performance profiles change.

# VI. Pillar 5: The LLM Prompt as a Formal Specification: A Critical Analysis

The ultimate success of the GDR protocol hinges on its most foundational premise: that a Large Language Model (LLM) prompt can function as an unambiguous, deterministic, and reliable specification for a task—effectively, a contract. This final section critically examines this assumption, drawing on research into prompt engineering, LLM determinism, and the inherent challenges of applying natural language processing to formal domains. The findings reveal significant technical and conceptual hurdles that challenge the viability of a naive implementation of the GDR concept.

### 6.1 State-of-the-Art in Prompt Engineering and Optimization

Prompt engineering has rapidly evolved from a niche art into a crucial discipline for controlling LLM behavior.<sup>73</sup> The field has established a range of techniques designed to improve the accuracy, consistency, and structure of model outputs. These include foundational methods like **zero-shot** and **few-shot prompting** (providing examples within the prompt), and more advanced strategies such as **chain-of-thought (CoT) prompting**, which guides the model to "think step-by-step" to improve performance on complex reasoning tasks.<sup>75</sup>

However, the very existence and complexity of this field underscores a fundamental problem: finding the "right" prompt is a difficult and often non-intuitive process. Minor variations in wording, formatting, or the choice of examples can lead to substantial differences in the model's response.<sup>77</sup> This has led to the rise of **Automatic Prompt Optimization (APO)**, a research area dedicated to systematically discovering effective prompts.<sup>79</sup> APO frameworks use a variety of optimization techniques, including:

- **Foundation Model-based Optimization:** Using an LLM to iteratively refine and improve prompts (meta-prompting).
- **Evolutionary Computing:** Applying genetic algorithms to mutate and evolve a population of candidate prompts.
- **Gradient-Based Optimization:** Treating discrete prompts as optimizable parameters in a differentiable system.
- **Reinforcement Learning:** Framing prompt editing as a series of actions in an RL environment, where the reward is based on the quality of the final output.

The fact that an entire subfield of AI research is dedicated to the complex, non-linear optimization problem of finding a good prompt is strong evidence against the assumption that a prompt can easily serve as a simple, unambiguous specification.

# 6.2 The Determinism Fallacy: Technical Barriers to Reproducibility at Temperature=0

A core requirement of the GDR protocol is that a given GDR, when processed by a specified LLM at temperature=0, will always regenerate the *identical* document. This property is essential for verification; if the output is not deterministic, it becomes impossible to simply hash the result and compare it to an expected value. However, a significant body of evidence from both industry practitioners and researchers demonstrates that setting temperature=0 does not guarantee deterministic outputs.<sup>83</sup>

This non-determinism is not a bug but an emergent property of the highly optimized, parallelized hardware and software stacks used to run modern LLMs. The key technical reasons include:

- 1. **Floating-Point Non-Associativity:** LLM inference involves billions of floating-point arithmetic operations. On modern GPUs, these operations are executed in parallel across thousands of cores. Floating-point addition is not perfectly associative (i.e., ) due to rounding errors. Because the order in which parallel threads complete and sum their results is non-deterministic, minuscule variations can arise in the final computed probabilities for the next token. If two tokens have extremely close probabilities, this tiny numerical "noise" can be enough to change their rank order, causing the greedy decoding process (argmax) to select a different token from one run to the next.<sup>83</sup>
- 2. **Mixture of Experts (MoE) Architectures:** Many state-of-the-art models, including those from the GPT-4 family (rumored) and open models like DeepSeek, use an MoE architecture. <sup>85</sup> In MoE models, a "gating network" routes each input token to one of a few specialized "expert" sub-models. For efficiency, inference is often batched, meaning requests from multiple users are processed simultaneously. The routing decision for a

token in one user's prompt can be influenced by the other tokens present in the same batch. This introduces a source of non-determinism at the level of an individual prompt, as its output can change depending on which other prompts were processed alongside it.<sup>85</sup>

3. **Hardware and Software Variations:** Different GPU models (e.g., A100 vs. H100), driver versions, or underlying deep learning library versions (like CUDA) can have slightly different implementations of mathematical operations, leading to small divergences in output even with identical inputs and model weights.<sup>83</sup>

This "determinism fallacy" represents a critical flaw in the naive conception of the GDR protocol. The core premise of regenerating an identical, bit-for-bit perfect document is technically fragile and likely unachievable with current and near-future LLM infrastructure. This necessitates a fundamental rethinking of the verification mechanism, moving from a simple hash comparison to a more nuanced measure of "semantic equivalence."

# 6.3 The "Prompt as Contract" Challenge: Ambiguity, Security, and Legal Precedent

Beyond the issue of determinism, treating a natural language prompt as a formal contract introduces a host of semantic, security, and legal challenges.

- Ambiguity: Natural language is inherently ambiguous, and LLMs often struggle to handle this. When faced with a vague or underspecified prompt, an LLM will typically make a plausible assumption and generate a response, rather than asking for clarification.<sup>78</sup> Research on prompt ambiguity shows that LLMs are highly sensitive to phrasing, and even small changes can lead to vastly different interpretations and outputs.<sup>77</sup> This makes a natural language prompt an unreliable foundation for a contract, where precision and lack of ambiguity are paramount.
- Security Vulnerabilities: LLM prompts are susceptible to prompt injection attacks, where malicious instructions embedded within what appears to be benign input can hijack the model's behavior, causing it to bypass safety guidelines, leak confidential information, or perform unintended actions.<sup>89</sup> If a GDR is a self-contained "job contract" that is executed autonomously, this vulnerability represents a critical security risk. An attacker could craft a GDR that appears to request a simple task but contains hidden instructions to corrupt the output or attack the executing agent's system.
- Legal and Contractual Complexity: The application of LLMs to formal legal documents
  has proven to be exceptionally challenging. LLMs often fail to grasp the nuanced,
  context-dependent language of contracts, leading to "confident mistakes" or
  hallucinations that are unacceptable in a legal setting. Furthermore, the legal
  framework surrounding LLM usage is unfavorable for such applications. The terms of

service for virtually all major LLM providers explicitly disclaim all warranties, state that the service is provided "AS-IS," and place the full responsibility and liability for the use of the output on the customer. 92 This makes it legally untenable to treat the output of an LLM, and by extension the prompt that generated it, as a binding agreement with any guarantee of correctness or fitness for a particular purpose.

These challenges suggest that the metaphor of a "prompt as a contract" is deeply flawed. A more accurate and robust architectural framing is to consider the prompt as a "blueprint." In this model, the legally and technically binding agreement is not the prompt itself, but the on-chain smart contract that governs its execution. The smart contract defines the parties, the reward, the deadline, and the objective verification criteria, while the GDR prompt serves as the set of instructions or specifications that are the *subject* of that contract. This separation of concerns isolates the robust, formal logic of the on-chain contract from the "soft," semantic, and potentially fallible nature of the natural language prompt.

# VII. Synthesis and Strategic Recommendations for the GDR Protocol

This comprehensive analysis of the research landscape across the five foundational pillars of the Generative Document Representation (GDR) concept reveals a project of significant ambition and novelty, positioned at the confluence of several major technological trends. The synthesis of these findings leads to a set of critical risks and strategic recommendations that should guide the future development of the GDR protocol.

### 7.1 GDR's Novelty and Positioning within the Research Landscape

The GDR concept's primary novelty does not lie in the invention of any single component technology, but in its unique synthesis and application of existing and emerging research.

- As a Compression Paradigm: GDR moves beyond the current focus of text compression (i.e., shortening context) to propose a new paradigm of *generative representation*, where a prompt serves as the canonical source for a complex document. This is a novel application for text, inspired by, but distinct from, the more mature field of generative visual compression.
- As a Communication Protocol: GDR is best positioned not as a competing transport or orchestration protocol, but as a high-level payload protocol. Its innovation lies in

- standardizing the *content* of agent communication—the semantic intent—which can then be carried by established protocols like ACP or A2A.
- As a DeAl Application: The GDR protocol provides a powerful, concrete use case that
  integrates the disparate components of the Decentralised Al ecosystem. It architects a
  "full stack" application that connects decentralized compute, autonomous economic
  agents, and DAO-based governance into a cohesive, value-generating system.

### 7.2 Key Technical Risks and Proposed Mitigation Strategies

The analysis has identified three primary technical risks that threaten the core assumptions of the GDR protocol. Addressing these risks proactively is essential for the project's success.

#### Risk 1: The Non-Determinism of LLM Inference

The foundational assumption of identical, bit-for-bit document regeneration at temperature=0 is technically unsound due to the nature of modern LLM infrastructure.83 Relying on this property for verification will lead to systemic failures.

- Mitigation Strategy: Pivot from "Identical Regeneration" to "Semantic Equivalence."
   The protocol's verification mechanism must be redesigned to be resilient to minor, non-deterministic variations in output. Instead of comparing hashes of the final document, verification should confirm that the generated output is semantically faithful to the creator's intent. This can be achieved by:
  - 1. Developing and standardizing robust **semantic equivalence metrics**, building upon concepts like the Semantic Reconstruction Effectiveness (SRE) proposed by Gilbert et al., which uses embedding similarity.<sup>21</sup>
  - 2. Incorporating an **LLM-as-a-judge** pattern, where a separate, trusted LLM evaluates the generated output against a set of criteria or a checklist of required attributes also included in the GDR packet.
  - 3. For structured data, verifying that the output conforms to a specified schema and that the extracted values are within acceptable bounds, rather than being character-for-character identical.

#### Risk 2: Prohibitive Overhead of Verifiable Computation

The computational cost of generating a full Zero-Knowledge proof for a large-scale LLM inference is currently too high to be economically viable for the vast majority of tasks.66 Mandating zkML for all GDR tasks would render the marketplace unusable.

Mitigation Strategy: Implement a "Tiered Verification" System. The GDR protocol
should allow the task creator to specify a required level of verification, tying the cost and
rigor of the proof to the economic value and security requirements of the task. This
creates a market-driven balance between trust and cost. A potential tiered system could
include:

- Tier O (Reputation-Based): No on-chain proof required. Payment is released based on the worker agent's established reputation within the network. Suitable for very low-value tasks.
- Tier 1 (Optimistic/Hardware-Based): Requires proof from a faster but less trustless method, such as a TEE attestation or an optimistic verification protocol with a short challenge window. Suitable for medium-value, latency-sensitive tasks.
- Tier 2 (Full Cryptographic Proof): Requires a full zk-SNARK proof of execution.
   Reserved for high-value, high-security, or privacy-sensitive tasks where the cost of proof generation is justified.

Risk 3: Ambiguity and Security of the "Prompt as Contract"

Treating a natural language prompt as a formal, binding contract is untenable due to the inherent ambiguity of language, the vulnerability to prompt injection attacks, and the lack of legal or technical guarantees from LLM providers.86

- Mitigation Strategy: Reframe the GDR as a "Blueprint" within a Smart Contract
  "Shell." The architecture must clearly separate the formal, on-chain agreement from the
  informal, off-chain instructions.
  - 1. The **Smart Contract** is the true, binding contract. It defines the immutable terms: the parties involved, the reward, the deadline, the hash of the specified LLM, and the objective conditions for payment (e.g., the successful on-chain verification of a Tier 2 proof).
  - 2. The GDR Prompt, stored on-chain or on a decentralized storage network like IPFS and referenced by the smart contract, is the blueprint or specification for the work to be done.
    - This separation of concerns makes the system robust. The smart contract provides the rigid, verifiable, and enforceable framework, while the prompt provides the flexible, semantic instructions for the generative task.

### 7.3 Strategic Opportunities and a Potential R&D Roadmap

By addressing these risks, the GDR protocol is well-positioned to become a standard for high-value, verifiable generative work in the rapidly emerging On-Chain AI Agent Economy. It provides a missing piece of the puzzle: a standardized way to package, trade, and verifiably execute complex, intent-driven tasks.

A potential R&D roadmap to realize this vision could be structured as follows:

- 1. Phase 1: Feasibility and Metric Development (Months 1-6)
  - **Objective:** Empirically quantify the non-determinism problem and develop robust evaluation metrics.
  - Key Activities:

- Conduct large-scale experiments across various open and closed-source LLMs (at temperature=0) to measure the frequency and magnitude of output variation for identical prompts.
- Develop and benchmark a suite of "semantic equivalence" metrics (e.g., embedding-based SRE, schema validation, LLM-as-a-judge frameworks).
- Publish findings to establish a clear, data-backed understanding of the determinism challenge.

#### 2. Phase 2: Protocol and Smart Contract Architecture (Months 4-9)

• **Objective:** Design the core components of the GDR protocol and its on-chain infrastructure.

#### Key Activities:

- Define the data structure for the "GDR Packet," including fields for the prompt, model identifier, verification tier, reward, and semantic equivalence criteria.
- Architect the suite of smart contracts for the job marketplace, including logic for posting tasks, escrowing funds, and handling automated payouts based on verification outcomes.
- Design the logic for the "Tiered Verification" system within the smart contracts.

#### 3. Phase 3: Proof-of-Concept and Integration (Months 7-18)

Objective: Build an end-to-end prototype demonstrating the integration of the DeAl stack.

#### Key Activities:

- Develop a basic worker agent (e.g., using the Fetch.ai AEA framework) capable of polling the smart contract marketplace for GDR tasks.
- Integrate with a decentralized compute provider (e.g., io.net) to allow the agent to provision GPU resources for LLM inference.
- Implement a simple Tier 0 (reputation) and Tier 1 (e.g., TEE-based, using a service like EigenCompute) verification flow.
- Integrate with an established agent communication protocol (e.g., ACP) to transport GDR packets.

#### 4. Phase 4: Advanced Verification and Security Hardening (Months 15-24)

• **Objective:** Implement high-security verification and harden the protocol against attacks.

#### Key Activities:

- Begin R&D on integrating a Tier 2 (zkML) verification option for high-value tasks, exploring partnerships with specialized companies (e.g., EigenCloud) or leveraging optimizing frameworks (e.g., ZKTorch).
- Conduct extensive security audits of the smart contracts and the GDR packet handling process to mitigate risks like prompt injection and economic exploits.
- Develop and deploy a DAO structure (e.g., using Aragon) to begin transitioning governance of the protocol to the community.

#### Works cited

- 1. What Are Al Agent Protocols? | IBM, accessed on October 10, 2025, https://www.ibm.com/think/topics/ai-agent-protocols
- 2. Agentic Al Communication Protocols: The Backbone of Autonomous ..., accessed on October 10, 2025, <a href="https://datasciencedojo.com/blog/agentic-ai-communication-protocols/">https://datasciencedojo.com/blog/agentic-ai-communication-protocols/</a>
- 3. Large Language Models: A Survey arXiv, accessed on October 10, 2025,
- https://arxiv.org/html/2402.06196v3
  Mitigating Visual Knowledge Forgetting in MLLM Instruction-tuning via
- Mitigating Visual Knowledge Forgetting in MLLM Instruction-tuning via Modality-decoupled Gradient Descent - arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2502.11740v1">https://arxiv.org/html/2502.11740v1</a>
- Mitigating Visual Knowledge Forgetting in MLLM Instruction-tuning via Modality-decoupled Gradient Descent - arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/pdf/2502.11740">https://arxiv.org/pdf/2502.11740</a>
- 6. Demystifying Reasoning Dynamics with Mutual Information: Thinking Tokens are Information Peaks in LLM Reasoning arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2506.02867v2">https://arxiv.org/html/2506.02867v2</a>
- 7. Generative Visual Compression: A Review arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2402.02140v1">https://arxiv.org/html/2402.02140v1</a>
- Unveiling the Future of Human and Machine ... Semantic Scholar, accessed on October 10, 2025, <a href="https://pdfs.semanticscholar.org/8874/cd2b166a1cc18642e0afd1ea44bbfefca8b1.pdf">https://pdfs.semanticscholar.org/8874/cd2b166a1cc18642e0afd1ea44bbfefca8b1.pdf</a>
- 9. [2402.02140] Generative Visual Compression: A Review arXiv, accessed on October 10, 2025, https://arxiv.org/abs/2402.02140
- 10. [Literature Review] Generative Visual Compression: A Review Moonlight, accessed on October 10, 2025, <a href="https://www.themoonlight.io/en/review/generative-visual-compression-a-review">https://www.themoonlight.io/en/review/generative-visual-compression-a-review</a>
- 11. Generative Latent Coding for Ultra-Low Bitrate Image Compression CVF Open Access, accessed on October 10, 2025, <a href="https://openaccess.thecvf.com/content/CVPR2024/papers/Jia\_Generative\_Latent\_Coding\_for\_Ultra-Low\_Bitrate\_Image\_Compression\_CVPR\_2024\_paper.pdf">https://openaccess.thecvf.com/content/CVPR2024/papers/Jia\_Generative\_Latent\_Coding\_for\_Ultra-Low\_Bitrate\_Image\_Compression\_CVPR\_2024\_paper.pdf</a>
- 12. Generative Image Compression by Estimating Gradients of the Rate-variable Feature Distribution arXiv, accessed on October 10, 2025, https://arxiv.org/html/2505.20984v1
- 13. High-Fidelity Generative Image Compression, accessed on October 10, 2025, <a href="https://proceedings.neurips.cc/paper\_files/paper/2020/hash/8a50bae297807da9e97722a0b3fd8f27-Abstract.html">https://proceedings.neurips.cc/paper\_files/paper/2020/hash/8a50bae297807da9e97722a0b3fd8f27-Abstract.html</a>
- 14. High-Fidelity Generative Image Compression, accessed on October 10, 2025, <a href="https://papers.neurips.cc/paper\_files/paper/2020/file/8a50bae297807da9e97722a0b3fd8f27-Paper.pdf">https://papers.neurips.cc/paper\_files/paper/2020/file/8a50bae297807da9e97722a0b3fd8f27-Paper.pdf</a>
- 15. Once-for-All: Controllable Generative Image Compression with Dynamic Granularity Adaptation | OpenReview, accessed on October 10, 2025, https://openreview.net/forum?id=z0hUsPhwUN
- 16. Semantics-Guided Generative Image Compression arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2505.24015v1">https://arxiv.org/html/2505.24015v1</a>

- 17. Semantic compression Wikipedia, accessed on October 10, 2025, <a href="https://en.wikipedia.org/wiki/Semantic compression">https://en.wikipedia.org/wiki/Semantic compression</a>
- 18. Semantic Compression with Large Language Models | Request PDF ResearchGate, accessed on October 10, 2025,
  <a href="https://www.researchgate.net/publication/377100201\_Semantic\_Compression\_with-Large-Language-Models">https://www.researchgate.net/publication/377100201\_Semantic\_Compression\_with-Large-Language-Models</a>
- 19. SCOPE: A Generative Approach for LLM Prompt Compression arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2508.15813v1">https://arxiv.org/html/2508.15813v1</a>
- 20. SCOPE: A Generative Approach for LLM Prompt Compression arXiv, accessed on October 10, 2025, <a href="https://www.arxiv.org/pdf/2508.15813">https://www.arxiv.org/pdf/2508.15813</a>
- 21. Semantic Compression With Large Language Models Computer Science, accessed on October 10, 2025, <a href="https://www.cs.wm.edu/~dcschmidt/PDF/Compression-with-LLMs-FLLM.pdf">https://www.cs.wm.edu/~dcschmidt/PDF/Compression-with-LLMs-FLLM.pdf</a>
- 22. Semantic Compression With Large Language Models arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/abs/2304.12512">https://arxiv.org/abs/2304.12512</a>
- 23. Semantic Compression With Large Language Models arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/pdf/2304.12512">https://arxiv.org/pdf/2304.12512</a>
- 24. Extending Context Window of Large Language Models via Semantic ..., accessed on October 10, 2025, <a href="https://aclanthology.org/2024.findings-acl.306/">https://aclanthology.org/2024.findings-acl.306/</a>
- 25. sarl/sarl-acl: FIPA Agent Communication Language for SARL GitHub, accessed on October 10, 2025, <a href="https://github.com/sarl/sarl-acl">https://github.com/sarl/sarl-acl</a>
- 26. Agent Communications Language Wikipedia, accessed on October 10, 2025, <a href="https://en.wikipedia.org/wiki/Agent\_Communications\_Language">https://en.wikipedia.org/wiki/Agent\_Communications\_Language</a>
- 27. Developing multi-agent systems with a FIPA-compliant agent framework, accessed on October 10, 2025, <a href="https://www.emse.fr/~boissier/enseignement/maop14/courses/readings/FIPA-JADE.pdf">https://www.emse.fr/~boissier/enseignement/maop14/courses/readings/FIPA-JADE.pdf</a>
- 28. Al Agent Communication from Internet Architecture Perspective: Challenges and Opportunities arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2509.02317v1">https://arxiv.org/html/2509.02317v1</a>
- 29. FIPA ACL Message Structure Specification FIPA.org, accessed on October 10, 2025, <a href="http://www.fipa.org/specs/fipa00061/SC00061G.html">http://www.fipa.org/specs/fipa00061/SC00061G.html</a>
- 30. FIPA ACL Message Structure Specification, accessed on October 10, 2025, <a href="http://www.fipa.org/specs/fipa00061/XC00061D.html">http://www.fipa.org/specs/fipa00061/XC00061D.html</a>
- 31. AGENT- COMMUNICATION LANGUAGES: IIIA-CSIC, accessed on October 10, 2025, <a href="https://www.iiia.csic.es/~puyol/SEIAD2001/publicacions/ACL-FIPA.doc.pdf">https://www.iiia.csic.es/~puyol/SEIAD2001/publicacions/ACL-FIPA.doc.pdf</a>
- 32. MCP vs A2A: A Guide to Al Agent Communication Protocols Auth0, accessed on October 10, 2025, <a href="https://auth0.com/blog/mcp-vs-a2a/">https://auth0.com/blog/mcp-vs-a2a/</a>
- 33. What is the Model Context Protocol (MCP)? | Cloudflare, accessed on October 10, 2025,
  - https://www.cloudflare.com/learning/ai/what-is-model-context-protocol-mcp/
- 34. What is the Model Context Protocol (MCP)? Model Context Protocol, accessed on October 10, 2025, <a href="https://modelcontextprotocol.io/">https://modelcontextprotocol.io/</a>
- 35. What is Agent Communication Protocol (ACP)? IBM, accessed on October 10, 2025, <a href="https://www.ibm.com/think/topics/agent-communication-protocol">https://www.ibm.com/think/topics/agent-communication-protocol</a>

- 36. The Agent Communication Protocol (ACP) and Interoperable Al Systems Macronet Services, accessed on October 10, 2025, <a href="https://macronetservices.com/agent-communication-protocol-acp-ai-interoperability/">https://macronetservices.com/agent-communication-protocol-acp-ai-interoperability/</a>
- 37. Agent Communication and Ontologies SmythOS, accessed on October 10, 2025, <a href="https://smythos.com/developers/agent-development/agent-communication-and-ontologies/">https://smythos.com/developers/agent-development/agent-communication-and-ontologies/</a>
- 38. Semantic Interoperability of Multi-Agent Systems in Autonomous Maritime Domains MDPI, accessed on October 10, 2025, <a href="https://www.mdpi.com/2079-9292/14/13/2630">https://www.mdpi.com/2079-9292/14/13/2630</a>
- 39. (PDF) A Review on Decentralized Artificial Intelligence in the Era of ..., accessed on October 10, 2025, <a href="https://www.researchgate.net/publication/380564678\_A\_Review\_on\_Decentralized">https://www.researchgate.net/publication/380564678\_A\_Review\_on\_Decentralized</a> d Artificial Intelligence in the Era of Large Models
- (PDF) SoK: Decentralized AI (DeAI) ResearchGate, accessed on October 10, 2025, <a href="https://www.researchgate.net/publication/386143556\_SoK\_Decentralized\_AI\_DeAI">https://www.researchgate.net/publication/386143556\_SoK\_Decentralized\_AI\_DeAI</a>
- 41. io.net | Decentralized GPU Ecosystem for Al Workloads Save Up to 70% io.net, accessed on October 10, 2025, <a href="https://io.net/">https://io.net/</a>
- 42. Top 10 Al Crypto Coins (2025): Projects Leading the Al Future, accessed on October 10, 2025, <a href="https://blog.aelf.com/posts/top-10-ai-crypto-coins">https://blog.aelf.com/posts/top-10-ai-crypto-coins</a>
- 43. Top Decentralized Compute Projects For Al Business Spheron's Blog, accessed on October 10, 2025, <a href="https://blog.spheron.network/top-decentralized-compute-projects-for-ai-business">https://blog.spheron.network/top-decentralized-compute-projects-for-ai-business</a>
- 44. SingularityNET Announces Global Winners of Over \$1M Deep Funding Grants for Developers of Beneficial AGI VKTR.com, accessed on October 10, 2025, <a href="https://www.vktr.com/the-wire/singularitynet-announces-global-winners-of-over-1m-deep-funding-grants-for-developers-of-beneficial-agi/">https://www.vktr.com/the-wire/singularitynet-announces-global-winners-of-over-1m-deep-funding-grants-for-developers-of-beneficial-agi/</a>
- 45. Top Web3 Al Projects | Onchain Magazine, accessed on October 10, 2025, <a href="https://onchain.org/magazine/top-web3-ai-projects/">https://onchain.org/magazine/top-web3-ai-projects/</a>
- 46. Al-Based Crypto Tokens: The Illusion of Decentralized Al? arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2505.07828v1">https://arxiv.org/html/2505.07828v1</a>
- 47. On-Chain Al Agent Economy: A Paradigm Shift for Web3 CV VC, accessed on October 10, 2025, <a href="https://www.cvvc.com/blogs/on-chain-ai-agent-economy-a-paradigm-shift-for-web3">https://www.cvvc.com/blogs/on-chain-ai-agent-economy-a-paradigm-shift-for-web3</a>
- 48. WEPIN Blog | How On-Chain Al Agents Are Transforming Web3 Wallets and Infrastructure, accessed on October 10, 2025, <a href="https://www.wepin.io/en/blog/onchain-ai-wallet-infra">https://www.wepin.io/en/blog/onchain-ai-wallet-infra</a>
- 49. Autonomous Al Agent Economies: Self-Governing Digital Entities Kava.io, accessed on October 10, 2025, <a href="https://www.kava.io/news/autonomous-ai-agent-economies-self-governing-digit">https://www.kava.io/news/autonomous-ai-agent-economies-self-governing-digit</a>

#### al-entities

- 50. Why Al Agents Need Web3 More Than Web3 Needs Al ... MEXC, accessed on October 10, 2025,
  - https://www.mexc.com/news/why-ai-agents-need-web3-more-than-web3-needs-ai-agents/122101
- 51. Autonomous Economic Agent (AEA) Meaning in Crypto | Tangem, accessed on October 10, 2025,
  - https://tangem.com/en/glossary/autonomous-economic-agent-aea/
- 52. Autonomous Economic Agent (AEA) Definition | CoinMarketCap, accessed on October 10, 2025.
  - https://coinmarketcap.com/academy/glossary/autonomous-economic-agent-aea
- 53. Composable Life: Speculation for Decentralized Al Life arXiv, accessed on October 10, 2025, https://arxiv.org/html/2508.20668v1
- 54. Decentralized autonomous organization Wikipedia, accessed on October 10, 2025, <a href="https://en.wikipedia.org/wiki/Decentralized autonomous organization">https://en.wikipedia.org/wiki/Decentralized autonomous organization</a>
- 55. DAOs Explained: Complete Guide to Decentralized Autonomous Organizations Rapid Innovation, accessed on October 10, 2025, <a href="https://www.rapidinnovation.io/post/daos-explained-ultimate-guide-to-decentralized-autonomous-organizations">https://www.rapidinnovation.io/post/daos-explained-ultimate-guide-to-decentralized-autonomous-organizations</a>
- 56. Transforming Business Through DAO and Al Integration TokenMinds, accessed on October 10, 2025, <a href="https://tokenminds.co/blog/ai/decentralized-ai">https://tokenminds.co/blog/ai/decentralized-ai</a>
- 57. Artificial Intelligence and Decentralized Autonomous Organizations: Where two worlds meet, accessed on October 10, 2025, <a href="https://wiprotechblogs.medium.com/artificial-intelligence-and-decentralized-autonomous-organizations-where-two-worlds-meet-2173312ae764">https://wiprotechblogs.medium.com/artificial-intelligence-and-decentralized-autonomous-organizations-where-two-worlds-meet-2173312ae764</a>
- 58. The Future of DAOs is Powered by AI | Aragon Resource Library, accessed on October 10, 2025, https://www.aragon.org/how-to/the-future-of-daos-is-powered-by-ai
- 59. A Gentle Introduction to zkML. What is "Zero-Knowledge Machine ..., accessed on October 10, 2025, https://opengradient.medium.com/a-gentle-introduction-to-zkml-8049a0e10a04
- 60. What is zkML? Explanation & Use Cases Datawallet, accessed on October 10, 2025, https://www.datawallet.com/crypto/what-is-zkml
- 61. Zero Knowledge Machine Learning | Axoncraze, accessed on October 10, 2025, https://www.axoncraze.com/zero-knowledge-machine-learning/
- 62. Why Verify? EZKL Blog, accessed on October 10, 2025, <a href="https://blog.ezkl.xvz/post/whyverify/">https://blog.ezkl.xvz/post/whyverify/</a>
- 63. State of Verifiable Inference & Future Directions Equilibrium Labs, accessed on October 10, 2025, <a href="https://equilibrium.co/writing/state-of-verifiable-inference">https://equilibrium.co/writing/state-of-verifiable-inference</a>
- 64. zkGPT: An Efficient Non-interactive Zero-knowledge Proof ... USENIX, accessed on October 10, 2025, https://www.usenix.org/system/files/usenixsecurity25-qu-zkgpt.pdf
- 65. vCNN: Verifiable Convolutional Neural Network Based on zk-SNARKs ResearchGate, accessed on October 10, 2025, <a href="https://www.researchgate.net/publication/377071773\_vCNN\_Verifiable\_Convolutio">https://www.researchgate.net/publication/377071773\_vCNN\_Verifiable\_Convolutio</a>

- nal Neural Network Based on zk-SNARKs
- 66. ZKTorch: Compiling ML Inference to Zero-Knowledge Proofs via Parallel Proof Accumulation arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2507.07031v1">https://arxiv.org/html/2507.07031v1</a>
- 67. [2410.22307] SVIP: Towards Verifiable Inference of Open-source Large Language Models, accessed on October 10, 2025, <a href="https://arxiv.org/abs/2410.22307">https://arxiv.org/abs/2410.22307</a>
- 68. EigenCloud Introduces Verifiable Al Solutions with EigenAl and EigenCompute MEXC, accessed on October 10, 2025, <a href="https://www.mexc.com/kk-KZ/news/eigencloud-introduces-verifiable-ai-solutions-with-eigenai-and-eigencompute/116712">https://www.mexc.com/kk-KZ/news/eigencloud-introduces-verifiable-ai-solutions-with-eigenai-and-eigencompute/116712</a>
- 69. EigenCloud Introduces Verifiable AI Solutions with EigenAI and ..., accessed on October 10, 2025, <a href="https://www.mexc.com/it-IT/news/eigencloud-introduces-verifiable-ai-solutions-with-eigenai-and-eigencompute/116712">https://www.mexc.com/it-IT/news/eigencloud-introduces-verifiable-ai-solutions-with-eigenai-and-eigencompute/116712</a>
- 70. Top Al Inference Providers Hyperbolic, accessed on October 10, 2025, <a href="https://hyperbolic.ai/blog/top-ai-inference-providers">https://hyperbolic.ai/blog/top-ai-inference-providers</a>
- 71. Decentralized LLM inference from your terminal, verified on-chain: r/selfhosted Reddit, accessed on October 10, 2025, <a href="https://www.reddit.com/r/selfhosted/comments/1m4ujre/decentralized\_llm\_inference\_from\_your\_terminal/">https://www.reddit.com/r/selfhosted/comments/1m4ujre/decentralized\_llm\_inference\_from\_your\_terminal/</a>
- 72. Decentralized LLM inference from your terminal, verified on-chain: r/LocalLLaMA Reddit, accessed on October 10, 2025, <a href="https://www.reddit.com/r/LocalLLaMA/comments/1m4u914/decentralized\_llm\_inference\_from\_your\_terminal/">https://www.reddit.com/r/LocalLLaMA/comments/1m4u914/decentralized\_llm\_inference\_from\_your\_terminal/</a>
- 73. A Systematic Survey of Prompt Engineering in Large Language Models: Techniques and Applications arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2402.07927v2">https://arxiv.org/html/2402.07927v2</a>
- 74. The Prompt Report: A Systematic Survey of Prompt Engineering Techniques arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/abs/2406.06608">https://arxiv.org/abs/2406.06608</a>
- 75. Prompt Engineering: Types and Optimizations Aussie AI, accessed on October 10, 2025, <a href="https://www.aussieai.com/research/prompt-engineering">https://www.aussieai.com/research/prompt-engineering</a>
- 76. What is Prompt Optimization? | IBM, accessed on October 10, 2025, https://www.ibm.com/think/topics/prompt-optimization
- 77. Answering Questions in Stages: Prompt Chaining for Contract QA arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2410.12840v1">https://arxiv.org/html/2410.12840v1</a>
- 78. How to Measure Prompt Ambiguity in LLMs Ghost, accessed on October 10, 2025,
  - https://latitude-blog.ghost.io/blog/how-to-measure-prompt-ambiguity-in-llms/
- 79. [2502.11560] A Survey of Automatic Prompt Engineering: An Optimization Perspective, accessed on October 10, 2025, <a href="https://arxiv.org/abs/2502.11560">https://arxiv.org/abs/2502.11560</a>
- 80. [2502.16923] A Systematic Survey of Automatic Prompt Optimization Techniques arXiv, accessed on October 10, 2025, https://arxiv.org/abs/2502.16923
- 81. A Survey of Automatic Prompt Engineering: An Optimization Perspective arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/html/2502.11560v1">https://arxiv.org/html/2502.11560v1</a>
- 82. Prompt optimization with two gradients for classification in large ..., accessed on

- October 10, 2025.
- https://pureadmin.qub.ac.uk/ws/files/651345892/Prompt\_optimization\_with\_two\_g radients for classification in large language models.pdf
- 83. Does Temperature 0 Guarantee Deterministic LLM Outputs ..., accessed on October 10, 2025, <a href="https://www.vincentschmalbach.com/does-temperature-0-guarantee-deterministic-llm-outputs/">https://www.vincentschmalbach.com/does-temperature-0-guarantee-deterministic-llm-outputs/</a>
- 84. Understanding why deterministic output from LLMs is nearly impossible Unstract, accessed on October 10, 2025, <a href="https://unstract.com/blog/understanding-why-deterministic-output-from-llms-is-nearly-impossible/">https://unstract.com/blog/understanding-why-deterministic-output-from-llms-is-nearly-impossible/</a>
- 85. [D] Non-deterministic behavior of LLMs when temperature is 0 : r/MachineLearning Reddit, accessed on October 10, 2025, <a href="https://www.reddit.com/r/MachineLearning/comments/lie15ev/d\_nondeterministiculor behavior of llms when/">https://www.reddit.com/r/MachineLearning/comments/lie15ev/d\_nondeterministiculor behavior of llms when/</a>
- 86. Can LLMs handle ambiguity in language? Milvus, accessed on October 10, 2025, <a href="https://milvus.io/ai-quick-reference/can-llms-handle-ambiguity-in-language">https://milvus.io/ai-quick-reference/can-llms-handle-ambiguity-in-language</a>
- 87. Teaching AI to Clarify: Handling Assumptions and Ambiguity in Language Models, accessed on October 10, 2025, <a href="https://shanechang.com/p/training-llms-smarter-clarifying-ambiguity-assumptions/">https://shanechang.com/p/training-llms-smarter-clarifying-ambiguity-assumptions/</a>
- 88. Do LLMs Understand Ambiguity in Text? A Case Study in Open-world Question Answering, accessed on October 10, 2025, <a href="https://arxiv.org/html/2411.12395v1">https://arxiv.org/html/2411.12395v1</a>
- 89. LLM Management Risks, Challenges And Opportunities for ..., accessed on October 10, 2025, <a href="https://medium.com/@senpubali7/navigating-the-llm-frontier-challenges-and-op-portunities-for-businesses-16681483c337">https://medium.com/@senpubali7/navigating-the-llm-frontier-challenges-and-op-portunities-for-businesses-16681483c337</a>
- 90. The Top Contract Al Challenges and How to Solve Them Artificial ..., accessed on October 10, 2025, <a href="https://www.artificiallawyer.com/2025/10/01/the-top-contract-ai-challenges-and-how-to-solve-them/">https://www.artificiallawyer.com/2025/10/01/the-top-contract-ai-challenges-and-how-to-solve-them/</a>
- 91. Legal Evalutions and Challenges of Large Language Models arXiv, accessed on October 10, 2025, <a href="https://arxiv.org/pdf/2411.10137?">https://arxiv.org/pdf/2411.10137?</a>
- 92. Navigating the LLM Contract Jungle: A Lawyer's Findings From an LLM Terms Audit, accessed on October 10, 2025, <a href="https://contractnerds.com/navigating-the-llm-contract-jungle-a-lawyers-findings-from-an-llm-terms-audit/">https://contractnerds.com/navigating-the-llm-contract-jungle-a-lawyers-findings-from-an-llm-terms-audit/</a>