

# The Grounding Gap: Why Agentic AI Fails and How Multimodal Reality Grounding Fixes It

Rajesh Iyer  
iyer70@gmail.com

## Abstract

Language implicates physical reality without understanding it—and unlike symbolic codes, language is combinatorially unbounded. This is the semantic explosion: enterprise AI has moved from symbols (policy codes, customer IDs) that are ungrounded but bounded, to language representations that are ungrounded and unbounded. “The warehouse near the state line,” “our Texas facility,” “that property we discussed”—infinite expressions can reference the same physical entity, with implicit semantics no schema captures. The grounding problem doesn’t merely persist from symbols to language; it explodes. We argue that multimodal reality grounding—resolving heterogeneous signals (structured, textual, visual, auditory, spatial) to canonical entities corresponding to physical reality—is the missing foundation layer for enterprise AI. This capability is distinct from traditional entity resolution: we identify three architectural requirements that existing approaches fail to meet: cross-modal projection (resolving audio, video, and documents to the same entity space as structured data), confidence-weighted governance (audit trails with provenance satisfying SR 11-7 and BCBS 239), and real-time agent consumption (sub-second resolution for autonomous decision-making). Without grounding, multi-agent systems exhibit a characteristic failure mode: each agent reasons over a different version of the world; each is locally rational; the system is globally incoherent. Through illustrative analyses in insurance and banking, we demonstrate how ungrounded representations propagate errors and how explicit grounding transforms agentic systems from demonstrations into deployable infrastructure. This is a position paper: we argue for a new architectural layer, not a new algorithm.

## 1. Introduction

Enterprise AI investments are underperforming. Despite billions spent on large language models, retrieval systems, and automation platforms, most organizations report that AI pilots fail to scale to production. The conventional explanations—model hallucination, insufficient training data, poor prompt engineering—miss

a deeper structural problem that becomes acute when multiple AI agents must coordinate on shared tasks.

### 1.1 From Bounded Symbols to Unbounded Language

Enterprise systems have evolved through two representation paradigms. Both fail to ground AI in physical reality—but they fail differently, and the second failure is worse.

**Symbolic Representations: Ungrounded but Bounded.** Traditional systems encode the world through rigid schemas—policy codes, customer IDs, counterparty identifiers, coordinates. These are symbols with no semantic content. TXK-001 doesn’t “know” it refers to a city straddling two states. CUST-77291 doesn’t “know” it represents the same corporation as LEI-549300XYZ in the trading system. They are pointers in databases, meaningful only to the systems that created them.

But symbolic systems have a saving grace: bounded vocabulary. TXK-001 lives in a finite code space with explicit schema. The universe of possible referents is enumerable. Cross-reference tables, however manually maintained, can in principle map every symbol to every other. The grounding problem exists, but its surface area is constrained.

**Language Representations: Ungrounded and Unbounded.** LLMs process the world through tokens and embeddings. The word “Texarkana” has meaning—but only statistical meaning derived from co-occurrence patterns in training data. “Meridian Holdings” embeds near “corporation” and “financial”—but the model does not understand that Meridian Holdings LLC, Meridian Holdings Inc., and Meridian Corp. are the same legal entity with one beneficial owner, one credit exposure, one AML risk profile.

Language implicates physical reality without understanding it—and unlike symbolic codes, language is combinatorially unbounded. “The warehouse near the state line,” “our Texas facility,” “that property we discussed last quarter,” “the building where the hail damage occurred”—infinite expressions can point at the same physical entity, with implicit semantics that no schema captures. The grounding problem doesn’t

merely persist from symbols to language; it explodes. The space of expressions requiring resolution scales with the expressive power of natural language itself.

This is the semantic explosion: moving from codes to language multiplies the surface area of what needs grounding by orders of magnitude, while providing no mechanism to perform that grounding.

**What’s Missing: Physical Reality.** The property at coordinates 33.4418, -94.0477 exists—it has a roof, a legal jurisdiction, a regulatory regime. The corporation “Meridian Holdings” exists—it has offices, employees, bank accounts, counterparty relationships. Symbolic systems point at these entities with bounded vocabularies. Language systems gesture at them with unbounded expressions. Neither resolves representations to physical entities with real-world consequences. Until AI systems can perform that resolution, they manipulate symbols, not meaning.

## 1.2 Why This Is Not Entity Resolution

A skeptical reader will ask: “Isn’t this just entity resolution with better marketing?” The question deserves a precise answer. Traditional entity resolution—matching records across databases using deterministic rules or probabilistic models—is a well-established discipline with mature vendor solutions. Our claim is not that entity resolution is insufficient, but that it addresses a different problem.

Entity resolution operates on structured records with defined schemas. It asks: “Is record A in system X the same entity as record B in system Y?” The inputs are rows in databases; the output is a match/no-match decision with optional confidence score.

Multimodal reality grounding operates on heterogeneous signals across modalities. It asks: “Do this call recording, this GPS coordinate, this video frame, this PDF excerpt, and this database record all refer to the same physical entity—and with what confidence, based on what evidence, suitable for what regulatory audit?” The inputs span audio spectrograms, image features, text embeddings, and structured codes; the output is a governed resolution with full provenance.

Three architectural requirements distinguish grounding from entity resolution:

1. **Cross-Modal Projection.** Traditional ER assumes structured inputs. Grounding must project unstructured signals—audio, video, handwritten documents—into the same resolution space as structured data. This requires modality-specific encoders (ASR for audio, vision models for images, OCR+layout for documents) that produce representations comparable across modalities.
2. **Confidence-Weighted Governance.** Traditional ER produces match decisions. Grounding must produce auditable assertions with provenance: which

**Table 1:** Entity resolution vs. multimodal grounding.

Dimension	Traditional ER	Grounding
Inputs	Structured records	All modalities
Temporal	Batch / eventual	Real-time
Output	Match decision	Governed assertion
Consumer	BI / analytics	Autonomous agents
Failure mode	Data inconsistency	World incoherence

signals contributed, with what extraction confidence, weighted by what regulatory salience. When a regulator asks “why did you assign Arkansas jurisdiction?”, the answer must trace to specific evidence, not “the model said so.”

3. **Real-Time Agent Consumption.** Traditional ER runs as batch or near-real-time processes feeding downstream analytics. Grounding must serve autonomous agents making sub-second decisions. An agent processing a claim cannot wait for overnight batch reconciliation; it needs grounded entity state now.

Table 1 summarizes these distinctions. The contrast is not incremental improvement but categorical difference: grounding addresses a problem that traditional ER was not designed to solve.

Existing structured resolution vendors (Informatica, Reltio) excel at matching records across databases but treat unstructured data as out of scope. AI vendors (OpenAI, Anthropic) provide embeddings and retrieval but not governed entity resolution with lineage. Digital twin platforms (NVIDIA Omniverse) model physical reality but don’t connect to enterprise master data. The grounding layer falls between vendor categories—which is precisely why it doesn’t exist.

## 1.3 The Agentic Failure Mode

The grounding problem becomes acute in multi-agent systems. Consider an enterprise claims handling system with five specialized agents: FNOL (first notice of loss), Photo Analysis, Policy Interpretation, Coverage Determination, and Correspondence Generation.

A naive objection: “This is just bad architecture. Share state properly and the problem disappears.” But the objection assumes structured, schema-aligned inputs. What state should be shared when the FNOL agent processes a phone call (“I’m at Texarkana... no, the Texas side... well, actually it straddles”), the Photo agent processes GPS metadata, the Policy agent processes a scanned PDF, and the Coverage agent applies jurisdictional rules?

The inputs are not rows in a database that can be joined on a key. They are heterogeneous signals—audio, images, documents, structured codes—that must be resolved to a common entity before any state can be

**Table 2:** Multi-agent failure mode—each agent interprets heterogeneous signals independently, producing globally incoherent outcomes.

Agent	Input Signal	Interpretation
FNOL	Call: "Texas side"	Jurisdiction: TX
Photo	GPS metadata	Location: AR
Policy	Scanned PDF	Endorsements: TX
Coverage	State rules DB	Limits: AR
Correspondence	All above	Wrong deductible

shared. Without grounding, "sharing state" means sharing inconsistent interpretations. Each agent reasons over a different version of the world. Each is locally rational; the system is globally incoherent.

This is why enterprises report "agents don't work—too risky." They're right, given the architecture. The problem isn't the agents. The problem is that agents cannot share a grounded world state when their inputs span modalities that no existing system resolves to common entities.

## 2. The Incarnation Problem

Real-world entities manifest through multiple "incarnations"—modality-specific representations with no inherent linkage. A single physical entity (a property, a corporation, a counterparty) projects into seven modality spaces, each with distinct extraction challenges and no authoritative cross-reference.

A single property or corporation may appear across seven modality categories, each with distinct extraction challenges:

**Structured Data.** Policy code TXK-001, ZIP 75501, FIPS 48037, coordinates 33.44,-94.04, reinsurance zone. Ten representations of one location, no authoritative linkage beyond manual cross-reference tables.

**Documents.** Scanned policy declarations, handwritten ACORD forms, contractor estimates. OCR extracts "Texarkana" (misspelled); layout parsing misses address fields in non-standard templates.

**Unstructured Text.** "The property near the state line," "our Texas facility" (actually AR), adjuster notes, email threads. Critical entity references that no structured system captures.

**Images.** Satellite tiles, claimant photos, drone captures. EXIF GPS may be present, absent, or incorrect. Semantic content (roof damage, building type) has no binding to policy systems.

**Audio.** Call recordings with pronunciation variation, ambient noise, conversational ambiguity. "I'm at Texarkana... no, the Texas side... well, actually it straddles."—meaning emerges from discourse, not transcription.

**Video.** Silent dashcam footage; security cameras with timestamps but no location; adjuster walkthroughs narrating "this side crosses into Arkansas." Multimodal gold—if extractable.

**Physical/IoT.** Sensors reporting to Building\_7A; digital twins in simulation environments; LIDAR point clouds. The grounding problem extends: how do you link atoms to bits?

## 3. Illustrative Analysis: Insurance

We present an illustrative analysis based on a composite scenario drawn from industry experience. The figures are estimates intended to demonstrate the mechanism of grounding failure, not measured outcomes from a specific deployment. Rigorous empirical validation remains future work.

**Scenario.** A hailstorm hits Texarkana—a city straddling Texas and Arkansas with distinct regulatory regimes. Over 72 hours, an insurer receives multiple incarnations of one property:

**Table 3:** Multiple incarnations of one property.

Modality	Incarnation
Structured	TXK-001, ZIP 75502, GPS coords
PDF	"4217 Oak St, Texarkana, TX"
Text	"warehouse near the state line"
Image	7 photos, EXIF GPS shows AR
Audio	"Texas side... well, straddles"
Video	"This side crosses into AR"

### 3.1 Without Grounding

**Day 1:** FNOL call transcribed as "Texarkana, Texas side." Representative searches policy system, finds 47 hits, selects one. Claim created under TX jurisdiction.

**Day 3:** Documents conflict—policy says TX, contractor estimate says AR, ACORD form says "Texarkana" (misspelled). Adjuster notes discrepancy; no escalation path; work continues.

**Day 5:** Field video uploaded: "This side crosses into Arkansas." Stored as email attachment. Never reviewed—no system flags video content for jurisdiction relevance.

**Day 14:** Arkansas regulator audit identifies wrong jurisdiction assignment. Reserve adjustment required; potential fine; rework costs. Estimated total impact: \$50K–\$90K per misgrounded claim of this complexity.

### 3.2 With Grounding

**Day 1:** Before any agent acts, grounding layer processes all available signals. GPS coordinates (highest confidence: 0.99) resolve to Miller County, AR via reverse

geocoding. Policy text indicates TX, but OCR extraction confidence (0.89) is lower than GPS. Audio transcript is ambiguous (0.67 confidence). Weighted aggregation: 73% AR, 27% TX—below automation threshold. System routes to human adjuster with evidence summary.

**Day 2:** Adjuster reviews evidence package, confirms AR jurisdiction. Resolution committed: Entity #4472 = Arkansas jurisdiction, 94% confidence post-human-verification. All downstream agents inherit grounded state.

**Day 6:** Claim closed. Correct jurisdiction. Zero rework. Zero regulatory risk.

**Table 4:** Illustrative claim processing comparison.

Metric	Without	With
Jurisdiction	Wrong	Correct
Reserve error	\$40–50K	\$0
Days to close	14+	6
Regulatory risk	High	None

## 4. Illustrative Analysis: Banking

The grounding problem manifests differently in banking, where entity fragmentation creates compliance risk rather than operational error.

### 4.1 The Incarnations of a Corporate Client

A corporate client exists across seven systems with no canonical linkage:

**Table 5:** Seven incarnations of one corporate client.

System	Incarnation
KYC/Onboarding	Meridian Holdings LLC
Core Banking	CUST-77291
Trading/Counterparty	LEI-549300XYZ
Lending	Meridian Corp.
Wire Transfer	Meridian Holdings Inc.
Call Center	Voice print #8821
Email Archive	“the Meridian account”

### 4.2 Without Grounding

**Day 1, 2:47 PM:** Fraud detection flags suspicious wire: \$2.3M to unfamiliar beneficiary from “Meridian Holdings Inc.” Alert assigned to AML analyst.

**Day 1, 4:15 PM:** AML analyst searches KYC system for “Meridian Holdings Inc.” No exact match. Finds “Meridian Holdings LLC”—assumes different entity. Requests enhanced due diligence on “new” client.

**Day 2, 9:30 AM:** Credit risk agent (automated) flags exposure concentration. Searches lending system for

“Meridian”—finds “Meridian Corp.” with \$45M facility. Treats as separate client. No aggregated exposure view.

**Day 3, 4:00 PM:** 72 hours elapsed. SAR filing deadline at risk. Investigation finally consolidated after compliance officer manually connects LLC, Inc., and Corp. Total exposure: \$78M across three “different” entities that were always one client.

### 4.3 With Grounding

**Day 1, 2:47 PM:** Wire flagged. Grounding layer resolves “Meridian Holdings Inc.” to canonical Entity #7291. System surfaces: also known as Meridian Holdings LLC (KYC), Meridian Corp. (Lending), LEI-549300XYZ (Trading). Total exposure: \$78M. Alert enriched automatically.

**Day 1, 6:00 PM:** SAR decision made with complete picture. Filed within 4 hours. No manual entity hunting.

**Table 6:** Illustrative investigation comparison.

Metric	Without	With
Time to resolution	72+ hours	4 hours
SAR deadline	At risk	Met
Analyst hours	60+	<10
Exposure accuracy	Fragmented	Complete

## 5. Reference Architecture

We propose a four-layer reference architecture for multi-modal reality grounding. This is an architectural pattern, not a product specification; implementations will vary based on organizational context, existing infrastructure, and regulatory requirements.

### 5.1 Ingestion Layer

Modality-specific encoders normalize heterogeneous signals into representations suitable for cross-modal resolution:

- **Structured data:** Schema normalization, code standardization
- **Documents:** OCR + layout parsing + confidence scoring
- **Text:** NER + entity linking + embedding generation
- **Images:** Feature extraction + EXIF parsing + geolocation
- **Audio:** ASR + diarization + entity extraction
- **Video:** Keyframe extraction + scene understanding + audio track processing

Each encoder outputs: (1) modality-specific representation, (2) extracted entity candidates, (3) extraction confidence score, (4) provenance metadata.

**Computational Considerations.** Sub-second resolution demands GPU-native execution. Running ASR, vision encoders, OCR, and embedding generation sequentially on CPU makes real-time agent consumption infeasible. Production implementations require: (1) parallel modality processing on GPU clusters, (2) incremental resolution that incorporates signals as they arrive rather than waiting for all modalities, (3) caching of entity embeddings and resolution decisions, and (4) tiered latency targets—synchronous for agent-blocking queries, asynchronous for enrichment. The architecture assumes infrastructure comparable to modern ML serving platforms (NVIDIA Triton, tensor-native data layers); CPU-bound serial processing is not viable for the latency requirements specified.

## 5.2 Resolution Layer

The resolution layer projects all representations into a shared space and determines entity identity. This is not nearest-neighbor embedding search—it is governed resolution with explicit confidence modeling.

**Confidence Model.** We require that confidence be composable, auditable, and modality-aware. The specific combination rule is implementation-dependent; we illustrate with a weighted Bayesian formulation. Each signal  $s_i$  contributes to the posterior probability of entity  $e$  with weight  $w_i$ :

$$P(e | S) \propto P(e) \prod_i P(s_i | e)^{w_i} \quad (1)$$

where weights encode three factors: extraction confidence (GPS: 0.99; faded OCR: 0.6), modality reliability (structured data more reliable than ASR in noisy environments), and regulatory salience (for jurisdiction, geolocation outweighs text mentions). The formalism is illustrative—we do not claim a uniquely correct combination rule. The contribution is the requirement that confidence be first-class, not the specific formula.

**Resolution Algorithm.** Given candidate entity set  $E$  and signal set  $S$ , compute  $P(e | S)$  for each entity  $e \in E$ . If  $\max P(e | S) > \theta_{\text{auto}}$  (automation threshold), commit resolution. If  $\theta_{\text{review}} < \max P(e | S) < \theta_{\text{auto}}$ , route to human review queue. If  $\max P(e | S) < \theta_{\text{review}}$ , flag as unresolvable pending additional signals.

**Worked Example.** Texarkana claim arrives with four signals: (1) FNOL transcript: “Texas side” (text, confidence 0.67—ambiguous phrasing); (2) Policy PDF: address “Texarkana, TX” (OCR, confidence 0.89); (3) Photo EXIF: GPS 33.4418, -94.0477 (structured, confidence 0.99); (4) Video transcript: “this side crosses into Arkansas” (ASR, confidence 0.82).

Resolution: GPS coordinates are highest-confidence signal; reverse geocoding returns Miller County, AR. Policy text indicates TX, but OCR confidence is lower than GPS. Video transcript corroborates AR. Weighted

aggregation:  $P(\text{AR}) = 0.73$ ,  $P(\text{TX}) = 0.27$ . Below  $\theta_{\text{auto}} = 0.85$ —route to human. Adjuster confirms AR. Resolution committed: Entity #4472 = Arkansas jurisdiction, confidence 0.94 post-verification.

## 5.3 Governance Layer

Every resolution is an assertion with full lineage:

- Source systems and extraction methods for each contributing signal
- Confidence scores at extraction and aggregation stages
- Human review decisions with reviewer identity and timestamp
- Downstream consumers that acted on the resolution

This lineage enables regulatory response: “Decision referenced Entity #4472, resolved from policy TXK-001 + GPS + video + human confirmation, 94% confidence.” Audit trail satisfies SR 11-7 (model risk management) and BCBS 239 (risk data aggregation) requirements.

**Human-in-the-Loop.** Low-confidence resolutions route to domain-specific queues: claims adjusters for jurisdiction disputes, KYC analysts for entity conflicts. Human decisions become labeled training signal, closing the feedback loop. Processing is asynchronous by default; synchronous for regulatory-critical paths where downstream action must wait for resolution.

## 5.4 Consumption Layer

APIs expose grounded entities to downstream systems:

- **GenAI/RAG:** Grounded retrieval—queries resolve to entities, not just documents
- **Agentic AI:** Entity-aware action—agents share grounded world state
- **Automation:** Self-healing resolution—confidence degradation triggers re-grounding

## 5.5 The Transformation

With grounding infrastructure, the five-agent claims system transforms:

**Table 7:** Multi-agent success with shared grounded state.

Agent	Input	Result
FNOL	Entity #4472	AR ✓
Photo	Entity #4472	AR ✓
Policy	Entity #4472	AR ✓
Coverage	Entity #4472	AR ✓
Correspondence	Entity #4472	Correct ✓

The agents didn’t get smarter. The world did. Every agent references the same grounded entity; coherence is architectural, not emergent.

## 5.6 Implementation Experience

The architecture described is not purely theoretical. We have implemented a proof-of-concept in **InsightGrid**, a GPU-native lakehouse built on Polars GPU, NVIDIA RAPIDS, and Ray. The platform leverages GPU Direct Storage (GDS) and RDMA to minimize data movement overhead, enabling sub-second resolution at scale.

InsightGrid embeds all media types—documents, images, audio, video—into a joint embedding space, then connects these representations to structured master data through ECM-like canonical indices. This unification of structured, semi-structured, and unstructured data within a single queryable architecture addresses the cross-modal projection requirement directly.

### Technical Stack:

- **Compute:** Polars GPU for DataFrame operations, NVIDIA RAPIDS for ML primitives, Ray for distributed orchestration
- **Storage:** Apache Iceberg table format, GPU Direct Storage for zero-copy data paths
- **Networking:** RDMA for GPU-to-GPU communication, bypassing CPU memory bottlenecks
- **Indexing:** Joint embedding space with ECM-like canonical indices linking embeddings to governed master data entities

Preliminary benchmarks demonstrate **30x+ performance improvements** over CPU-based implementations for typical resolution workloads. The performance gains derive from redesigned data paths—not simply “adding GPUs” to existing architectures, but rethinking how data flows from storage through compute to resolution. Detailed empirical validation is ongoing.

## 6. Implementation Pathway

For enterprise practitioners evaluating grounding infrastructure, we outline a phased implementation approach. This is not a product pitch—it is a practical roadmap based on implementation experience.

### 6.1 Phase 1: Foundation (Weeks 1–6)

**Objective:** Establish grounding infrastructure for structured data and documents—the highest-value, lowest-risk starting point.

#### Activities:

- Deploy GPU-native lakehouse (InsightGrid or equivalent) with Iceberg table format
- Integrate existing structured resolution investments (MDM, ER outputs) as seed entity graph
- Implement document ingestion pipeline: OCR + layout parsing + confidence scoring
- Connect 2–3 source systems (e.g., policy admin, claims, documents)

- Establish canonical entity index with governance metadata

**Deliverable:** Grounded entity resolution for structured + document modalities, serving batch and near-real-time queries.

**Success Criteria:** Resolution latency <500ms for 95th percentile; governance lineage complete for all resolutions; human review queue operational.

### 6.2 Phase 2: Multimodal Extension (Weeks 7–12)

**Objective:** Extend grounding to audio and image modalities; enable real-time agent consumption.

#### Activities:

- Deploy ASR pipeline for call recordings with entity extraction
- Implement image ingestion: EXIF parsing, geolocation, feature extraction
- Tune confidence weights by modality based on Phase 1 human review patterns
- Build real-time resolution API for agent consumption (<100ms target)
- Integrate with one agentic workflow (e.g., FNOL processing)

**Deliverable:** Four-modality grounding (structured, documents, audio, images) with real-time API.

**Success Criteria:** Agent workflow demonstrates coherent multi-modal resolution; error rate reduction measurable against baseline.

### 6.3 Phase 3: Production Hardening (Weeks 13–18)

**Objective:** Scale to production volumes; establish operational governance.

#### Activities:

- Performance optimization for production throughput
- Implement confidence degradation monitoring and re-grounding triggers
- Build regulatory reporting dashboards (SR 11-7, BCBS 239 compliance)
- Extend to video modality if high-value use cases identified
- Document operational runbooks and escalation procedures

**Deliverable:** Production-grade grounding infrastructure with full governance.

**Success Criteria:** SLA compliance at production scale; regulatory audit readiness demonstrated; ROI metrics validated.

### 6.4 Investment Framework

**ROI Framing:** A mid-sized P&C carrier processing 500K claims annually with 3% jurisdictional ambiguity rate

**Table 8:** Illustrative investment framework by phase.

Phase	Duration	Team	Infrastructure
Foundation	6 weeks	4–6 FTE	GPU cluster (dev)
Extension	6 weeks	6–8 FTE	GPU cluster (staging)
Hardening	6 weeks	8–10 FTE	Production infra

faces \$7.5M–\$13.5M annual exposure from misgrounding (at \$50K–\$90K per complex misgrounded claim). A bank with 10,000 corporate clients fragmented across 5+ systems averages 15 hours analyst time per AML investigation due to entity hunting; grounding reduces this to <2 hours.

The payback calculation is domain-specific, but the pattern holds: grounding converts analyst hours into automated resolution, and converts regulatory risk into auditable governance.

## 6.5 Build vs. Buy Landscape

No vendor currently offers complete multimodal reality grounding. The market gap:

**Table 9:** Competitive landscape for grounding capabilities.

Vendor	Strength	Gap
Informatica, Reltio	Structured ER	Cross-modal, real-time
Palantir Foundry	Governed ontology	Auto multimodal resolution
OpenAI, Anthropic	Embeddings, retrieval	Governance, lineage
NVIDIA Omniverse	Digital twin physics	Enterprise MDM

The implementation pathway assumes organizations will **build** the grounding layer, potentially using InsightGrid or similar GPU-native infrastructure, while **integrating** existing investments in structured resolution and AI embeddings. This is infrastructure work, not vendor selection.

## 7. Limitations and Open Problems

We are candid about limitations. Multimodal reality grounding is not a panacea—it focuses human judgment on genuinely ambiguous cases rather than eliminating human involvement entirely. Several hard problems remain open:

**Low-Confidence Deadlock.** When resolution confidence sits between thresholds (e.g., 55%—too low for automation, too high for clear rejection), the system can stall. Current mitigation: time-boxed escalation with “best available” fallback and explicit uncertainty propagation to downstream agents. Better solutions may re-

quire active information gathering—querying additional sources or requesting specific evidence from users.

**Modality Conflict.** When GPS says Arkansas and notarized policy says Texas, which wins? Regulatory salience weighting provides a framework, but edge cases require domain-specific arbitration rules that must be maintained alongside grounding infrastructure. The meta-problem: who governs the governance rules?

**Adversarial Entities.** Shell companies with intentionally similar names. Properties undergoing jurisdictional rezoning. Merged corporations with legacy identifiers. The canonical entity itself may be contested or in flux. Grounding cannot resolve what reality has not yet settled—and sophisticated actors may exploit this ambiguity.

**Cold Start.** Grounding quality depends on accumulated resolutions and feedback. New entities, new modalities, and new domains start with sparse linkage. The architecture assumes bootstrapping from existing structured resolution investments, which may themselves be incomplete or inconsistent. Phase 1 of the implementation pathway addresses this by seeding the entity graph from existing MDM/ER outputs.

**Empirical Validation.** This paper presents architecture and illustrative analysis, not measured outcomes. Rigorous validation requires controlled deployment with before/after metrics, statistical significance testing, and longitudinal assessment of grounding quality over time. We identify this as essential future work.

**Value Proposition.** The business case for grounding infrastructure depends on claim complexity distribution. If 80% of claims are unambiguous (single jurisdiction, single entity, clear documentation), grounding automates only the remaining 20%—valuable, but not transformative. If ambiguity is pervasive, the value proposition strengthens. Empirical measurement of ambiguity rates across domains is needed.

## 8. Implications

**For Enterprise Architecture.** If the argument holds, grounding is a prerequisite layer for agentic AI—not an optimization to add later. Organizations deploying multi-agent systems without reality grounding will encounter the failure modes described: locally rational agents producing globally incoherent outcomes. The grounding layer should be architectural infrastructure, not an afterthought.

**For AI Governance.** Regulators increasingly require explainability. Without grounding lineage, explanations are “the model retrieved documents” or “the agent decided.” With grounding: “Decision referenced Entity #4472, resolved from policy TXK-001 + GPS + video, 94% confidence, human-confirmed.” This is auditable.

**For Practitioners.** The reference architecture is imple-

mentable today. InsightGrid demonstrates the approach: joint embedding of all modalities, indexed against governed canonical entities, with sub-second resolution serving agentic workloads. The performance economics of GPU acceleration make real-time grounding feasible where CPU-based approaches could not meet latency requirements.

**For Vendors.** Multimodal reality grounding falls between existing vendor categories. Structured resolution vendors could extend into unstructured modalities. AI vendors could add governance and lineage. Digital twin vendors could connect to enterprise master data. The opportunity exists for a new category—or for existing vendors to expand scope.

**For Research.** Open problems include: confidence modeling for heterogeneous signals, active information gathering for low-confidence resolution, adversarial robustness against entity manipulation, and empirical measurement of grounding quality over time. We hope this paper motivates further work.

## 9. Conclusion

Enterprise AI has progressed from data representations to language representations—but neither grounds AI in physical reality. Codes are symbols without semantics. Tokens are statistics without referents. The property in Texarkana exists, with a roof that hail can damage, a jurisdiction that determines liability. The corporation Meridian Holdings exists, with a single credit exposure across seven system incarnations. Until AI connects its representations to that reality, it manipulates symbols, not meaning.

We have argued that multimodal reality grounding—resolving heterogeneous signals to canonical, governed entities corresponding to physical reality—is the missing foundation layer for enterprise AI. This is distinct from traditional entity resolution: cross-modal projection, confidence-weighted governance, and real-time agent consumption constitute architectural requirements that existing approaches do not address.

The case is not proven—empirical validation is essential future work. But the architecture is implementable: InsightGrid demonstrates that GPU-native infrastructure can achieve the performance requirements for real-time agent consumption while maintaining governance lineage. Organizations that build grounding infrastructure transform what requires armies of analysts into automated, auditable, governable AI pipelines. The ones that don't will keep asking why their agents can't agree on jurisdiction, or why three agents investigated three "different" clients who were always one.

Agents don't need more autonomy—they need a shared reality.

## References

- [1] Harnad, S. (1990). The Symbol Grounding Problem. *Physica D*, 42, 335–346.
- [2] Bisk, Y. et al. (2020). Experience Grounds Language. *EMNLP 2020*.
- [3] Chen, Y. et al. (2023). A Survey on Multimodal Knowledge Graphs: Construction, Completion and Applications. *Mathematics*, 11(8).
- [4] Zeakis, A. et al. (2023). Pre-trained Embeddings for Entity Resolution: An Experimental Analysis. *PVLDB*, 16(9), 2225–2238.
- [5] Mudgal, S. et al. (2018). Deep Learning for Entity Matching: A Design Space Exploration. *SIGMOD 2018*.
- [6] Bordes, A. et al. (2013). Translating Embeddings for Modeling Multi-relational Data. *NeurIPS 2013*.
- [7] NVIDIA (2024). Omniverse Platform Documentation.
- [8] Palantir Technologies (2024). Foundry Ontology: Building the Operational Data Layer. Technical Documentation.
- [9] Board of Governors of the Federal Reserve System (2011). SR 11-7: Guidance on Model Risk Management.
- [10] Basel Committee on Banking Supervision (2013). BCBS 239: Principles for Effective Risk Data Aggregation and Risk Reporting.
- [11] Liu, Y. et al. (2019). MMKG: Multi-Modal Knowledge Graphs. *ESWC 2019*.
- [12] LeCun, Y. (2022). A Path Towards Autonomous Machine Intelligence. *OpenReview*.