

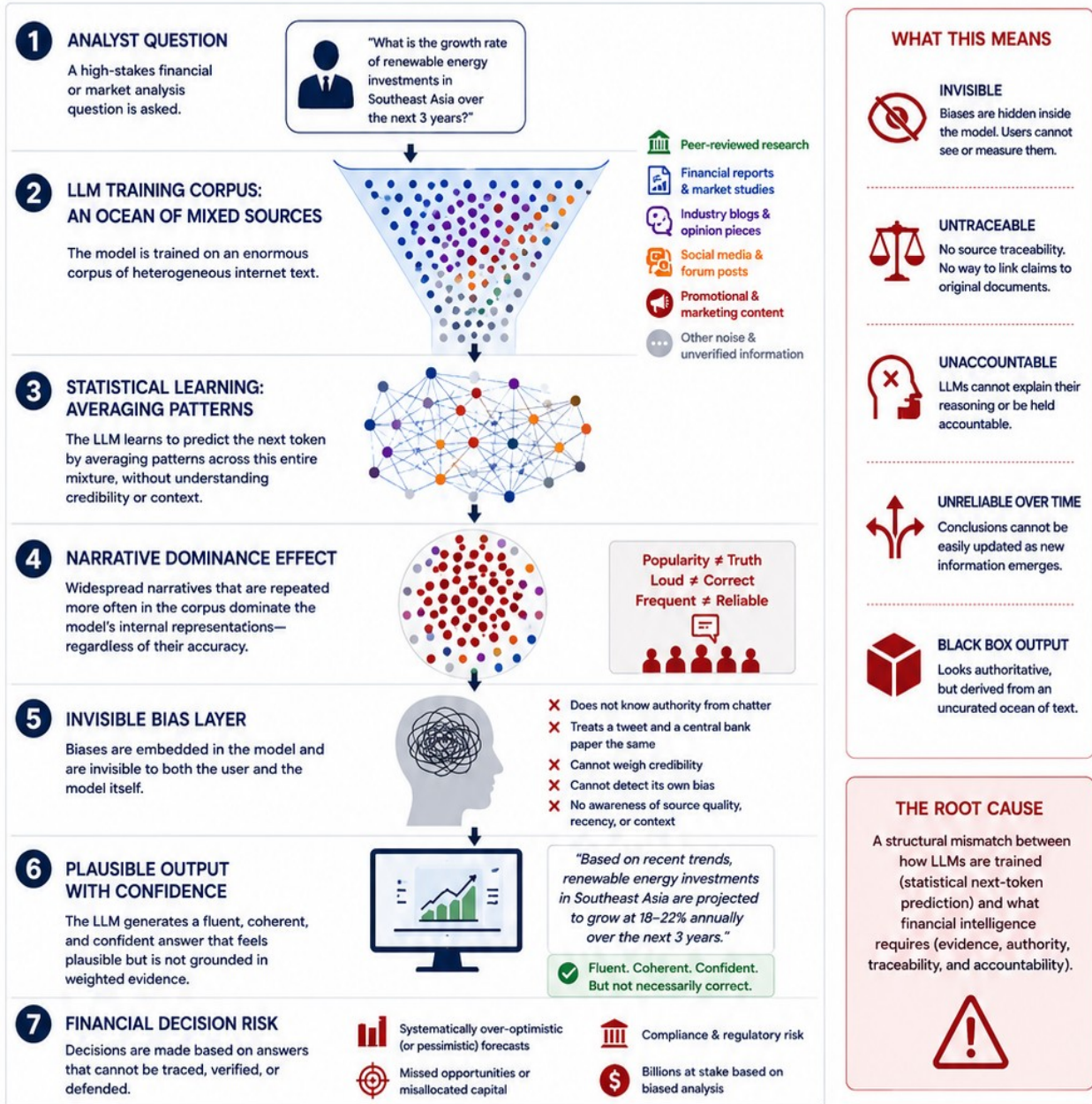
The Augmented Analyst: An AI Exoskeleton for Finance



THE PROBLEM: INVISIBLE BIASES IN LLMs FOR FINANCIAL AND MARKET ANALYSIS

LLMs learn from the statistical prevalence of text,
not the truth or the weight of evidence.

The result: plausible but biased answers with no traceability.

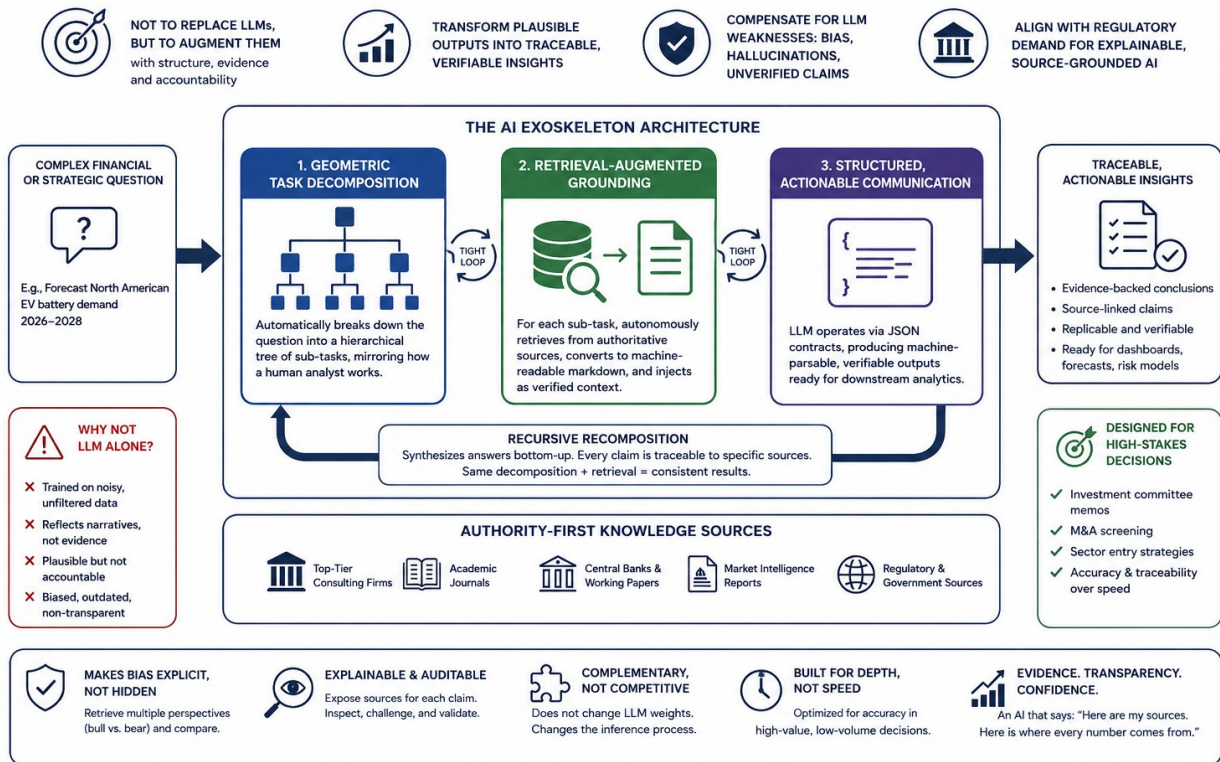


THE SOLUTION IS NOT TO DISCARD LLMs.
It is to surround them with a layer that enforces grounding, enables traceability, and imposes a structured decomposition and synthesis process.

THIS IS THE MISSION OF THE AI EXOSKELETON.

AI EXOSKELETON VISION

Evidence-Based Augmented Intelligence for Financial Decisions



Vision

The vision of the AI Exoskeleton is not to replace large language models nor to simply wrap them with a nicer interface. Instead, it is to structurally augment them by building an active, reasoning orchestration layer that systematically compensates for the fundamental weaknesses of pure LLMs while preserving their generative power. In the domain of financial analysis, market forecasting, benchmarking, and investment screening, these weaknesses are not marginal. LLMs trained on broad, unfiltered internet data inevitably internalize barroom chatter, journalistic oversimplifications, and socially prevalent but empirically fragile beliefs. When asked to forecast a sector’s growth, value a stock, or identify emerging markets, a standalone LLM produces outputs that are plausible but not accountable. Its answers are emergent averages of whatever narratives dominate its training corpus, not reasoned conclusions grounded in verifiable authoritative sources.

The AI Exoskeleton changes this by introducing a hybrid cognitive architecture composed of three tightly integrated innovations. The first is geometric task decomposition, meaning that any complex financial or strategic question is automatically broken down into a hierarchical structure of sub-tasks. For example, forecasting North American electric vehicle battery demand from 2026 to 2028 would decompose into regional policy analysis, raw

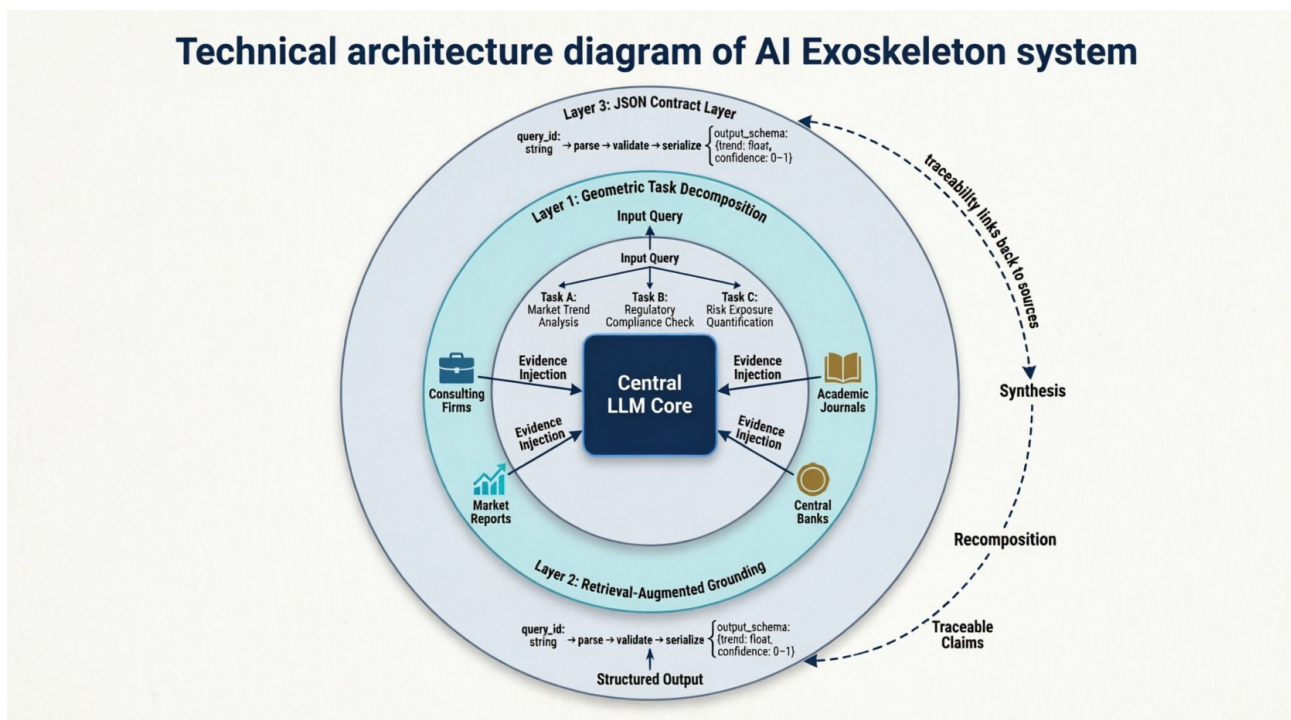
material supply forecasts, competitor capacity expansion plans, and so on. This mirrors how a diligent analyst would work, not how a stochastic parrot would guess. The second innovation is retrieval-augmented grounding, where for each sub-task the system autonomously retrieves documentation from pre-defined authoritative sources such as top-tier consulting firms, academic journals, central bank working papers, and recognised market intelligence reports. The retrieved content is converted to machine-readable markdown and injected as context, so that the LLM never answers from memory alone but on the basis of evidence that the orchestrator has explicitly fetched and validated at source. The third innovation is structured, actionable communication, enforced via JSON contracts between the orchestrator and the LLM. This turns the LLM from a conversational engine into a planning and execution component whose outputs are machine-parsable, verifiable, and immediately reusable in downstream analytics such as forecasting models, scoring dashboards, or risk matrices.

What makes this genuinely innovative, and not just another retrieval-augmented generation pipeline, is the tight loop between decomposition, retrieval, and recomposition. In conventional RAG, the LLM is given a few retrieved documents and asked a single, often broad question. The Exoskeleton, by contrast, recursively decomposes the original ask, retrieves highly specific documents for each leaf sub-task, synthesizes answers bottom-up, and finally produces a conclusion that is traceable because every claim links back to a specific retrieved source, and replicable because the same decomposition and retrieval steps yield consistent outputs even with non-deterministic LLMs.

Addressing the sceptic's concerns directly, one might ask whether LLMs are already improving and whether fine-tuning could solve the bias problem. Fine-tuning on financial data reduces some surface errors but does not eliminate deep structural biases, because those biases are embedded in the very distribution of the training data. The Exoskeleton does not rely on changing the LLM's weights; it changes the inference-time process, which is a complementary and provably different approach. Another concern is that authoritative documents can also be biased or outdated. That is true, but the Exoskeleton is architected to expose sources, not to sanctify them. The investor or analyst can inspect which documents were used for each claim, and the system can also be configured to retrieve multiple competing sources, for instance bull versus bear reports, and explicitly contrast them, turning bias from a hidden defect into an explicit risk variable. A third objection is that this approach adds latency and complexity. The value proposition, however, is explicitly for high-stakes, low-volume decisions such as quarterly investment committee memos, mergers and acquisitions screening, or sector entry strategies, where accuracy and

traceability matter more than sub-second response times. For real-time trading, other tools exist; the Exoskeleton is a decision support system, not a high-frequency trading bot.

The ultimate vision is an AI that does not hallucinate confidently. Instead, it says: I searched the following four authoritative reports; based on them, here is my analysis, and here is exactly where each number comes from. In a world where financial professionals are asked to do more with less, and where regulatory pressure for explainable AI is rising through the EU AI Act and proposed SEC rules, the AI Exoskeleton offers a pragmatic, defensible, and immediately implementable path toward truly evidence-based augmented intelligence. This white paper describes the architecture, validates its outputs against standard LLM baselines, and outlines concrete use cases in financial forecasting, market benchmarking, stock analysis, and investment area identification. The reader is invited to see the Exoskeleton not as a futuristic vision, but as a working prototype ready for deployment in professional financial environments.



The Problem: Invisible Biases in LLMs for Financial and Market Analysis

Large language models have demonstrated remarkable fluency and reasoning capabilities across a wide range of domains. Their strength, however, is also their fundamental weakness when applied to high-stakes financial and market analysis. An LLM is trained on an enormous corpus of text drawn from the public internet, which includes not only peer-reviewed research and reputable financial reports but also unverified blog posts, forum debates, promotional content, and casual conversations. The model learns to predict the next token by averaging patterns across this entire heterogeneous mixture. As a result, its outputs reflect not the truth or the weight of evidence, but rather the statistical prevalence of certain claims and narratives within its training data. If a particular belief is widespread on the internet regardless of its factual accuracy, the LLM will reproduce it as if it were a reasonable conclusion. This is not a bug that can be fixed by simple fine-tuning; it is a structural consequence of how language models are trained.

In the context of financial forecasting, market benchmarking, and investment analysis, this problem becomes acute. Consider a scenario where an analyst asks a standard LLM to project the growth rate of renewable energy investments in Southeast Asia over the next three years. The model will draw upon news articles, promotional materials from industry associations, social media discussions, and a handful of genuine market studies. The sheer volume of optimistic, headline-driven content may drown out more cautious or academically rigorous analyses. The LLM will produce a forecast that feels plausible and internally consistent, but it will be unable to cite which sources informed which numbers, nor will it have any mechanism to prioritize authoritative studies over popular opinion. The result is a prediction that may be systematically over-optimistic or, in other cases, overly influenced by recent but unrepresentative events. The problem is further compounded by the fact that these biases are invisible. Unlike a human analyst who can consciously weigh the credibility of a source, an LLM has no such awareness. It does not know that a tweet carries less evidential weight than a central bank working paper; it treats both as strings of tokens to be statistically modeled.

Another dimension of the issue is that biases embedded in widely shared narratives are often culturally and socially entrenched. For example, certain industries or asset classes may be subject to persistent myths about their risk-return profile, or geographical markets may be stereotyped based on outdated or incomplete information. Because these narratives appear repeatedly across the training corpus, the LLM internalizes them as

default assumptions. When asked to perform a market benchmark or to identify promising investment areas, the model will inadvertently reproduce these stereotypes, potentially leading to missed opportunities or undue risk exposure. A sceptical reader might argue that a human analyst is also subject to biases, and that is true. However, the human analyst can be trained to recognize and correct for known cognitive biases, and can be held accountable for their reasoning. An LLM offers no such transparency. Its biases are emergent properties of its training data, and they cannot be easily traced or corrected after the fact.

Furthermore, the problem is not solved by increasing the size of the model or the volume of training data. Larger models become even more adept at modeling statistical regularities, but they do not acquire an inherent ability to distinguish authority from chatter. In fact, larger models can become more confident in their outputs, making their errors harder to detect. Fine-tuning on financial domain data can help align the model with preferred styles or formats, but it does not eliminate the foundational issue: the model still has no native mechanism to retrieve, prioritize, or ground its answers in externally validated documents at inference time. It relies entirely on what it already memorized during training, which includes both the signal and the noise.

For professional financial applications, this is unacceptable. Investment decisions, corporate strategy, and regulatory compliance require not just plausible answers but demonstrably defensible ones. The analyst or portfolio manager must be able to trace each claim back to a source, to compare competing evidence, and to update conclusions when new information arrives. A pure LLM offers none of these capabilities. It provides a black box that outputs text that looks authoritative but is statistically derived from an uncurated ocean of text. The risk is not merely that the model will occasionally be wrong; it is that its errors will be systematically biased toward popular but flawed narratives, and that these errors will be presented with the same fluency and confidence as correct information. In a field where billions of dollars hinge on the quality of analysis, relying on such a tool without structural augmentation is a form of negligence.

Therefore, the problem that the AI Exoskeleton addresses is not a marginal improvement in accuracy or a reduction in hallucinations. It is the fundamental incompatibility between the way LLMs are trained and the requirements of evidence-based financial intelligence. The solution is not to discard LLMs, which remain extraordinarily capable generators of language and reasoning templates. The solution is to surround them with a layer that forces grounding, enables traceability, and imposes a structured decomposition and

synthesis process. Only then can an LLM be transformed from a statistical mimic of internet chatter into a reliable partner for high-consequence financial analysis.

The AI Exoskeleton: Architecture of a Hybrid Financial Agent

The AI Exoskeleton is not a single model nor a simple software wrapper. It is an orchestration layer that sits between the user and one or more large language models, fundamentally altering how the LLM is invoked and how its outputs are produced. The architecture rests on four logical components that work in a closed loop: problem decomposition, retrieval and grounding, structured generation, and synthesis. Each component is designed to compensate for a specific weakness of standalone LLMs, and together they form a cognitive exoskeleton that augments rather than replaces the underlying language model.

The process begins when a user submits a complex financial or strategic question. This could be a request to forecast the revenue growth of a specific automotive supplier over the next five years, or to benchmark the competitive landscape of telemedicine platforms in Southeast Asia, or to identify underappreciated investment opportunities in the European renewable energy sector. Instead of sending this request directly to an LLM, the Exoskeleton first invokes a dedicated decomposition module. This module prompts the LLM to break the original question into a hierarchical tree of sub-tasks. For a revenue forecast, the root task might decompose into sub-tasks such as historical revenue analysis, addressable market sizing, regulatory environment assessment, and competitor intelligence. Each of those can be further decomposed. The decomposition is guided by a structured JSON contract, ensuring that the output is machine-parsable and that the tree is both comprehensive and non-redundant. The depth and breadth of the decomposition can be configured by the user, allowing a trade-off between analytical thoroughness and response latency.

Once the decomposition tree is complete, the Exoskeleton moves to the retrieval and grounding phase. For each leaf sub-task, the system constructs a search query derived from the sub-task description and from metadata such as the domain and subdomain of the original question. This query is then executed against a curated set of authoritative sources. In the current implementation, these sources include publications from major consulting firms, academic databases, central bank working papers, and recognized market intelligence providers. The Exoskeleton does not rely on the LLM to recall or infer these documents; it actively downloads them, typically in PDF format, from publicly accessible or

subscribed repositories. Each downloaded document is then converted into plain markdown text using either a local converter or a remote optical character recognition service. The resulting text is stored locally for future reuse, avoiding repeated downloads and conversions of the same document. Importantly, the system maintains a hash-based index to detect duplicate documents across different URLs and to track which documents have already been processed. For each leaf sub-task, the Exoskeleton retrieves a set of relevant documents, but it does not simply concatenate them. Instead, it selects the most promising documents based on query relevance and, if desired, can also retrieve competing sources that offer contrasting viewpoints.

With the decomposition tree and the retrieved documents in hand, the Exoskeleton enters the structured generation phase. For each leaf sub-task, the system constructs a prompt that includes the sub-task description, the relevant portions of the retrieved documents, and a strict instruction to respond in a predefined JSON format. The prompt is sent to the LLM, which is configured with a temperature of zero to minimize randomness. The LLM's response is a JSON object containing the answer to the sub-task, along with citations linking specific claims back to the source documents. Because the prompt explicitly forbids the LLM from relying on its own internal knowledge alone, the model is forced to ground its answer in the provided documents. If the documents do not contain sufficient information to answer the sub-task, the LLM is instructed to state that gap explicitly rather than to hallucinate. This is a critical departure from standard LLM usage, where the model is implicitly encouraged to produce an answer no matter what.

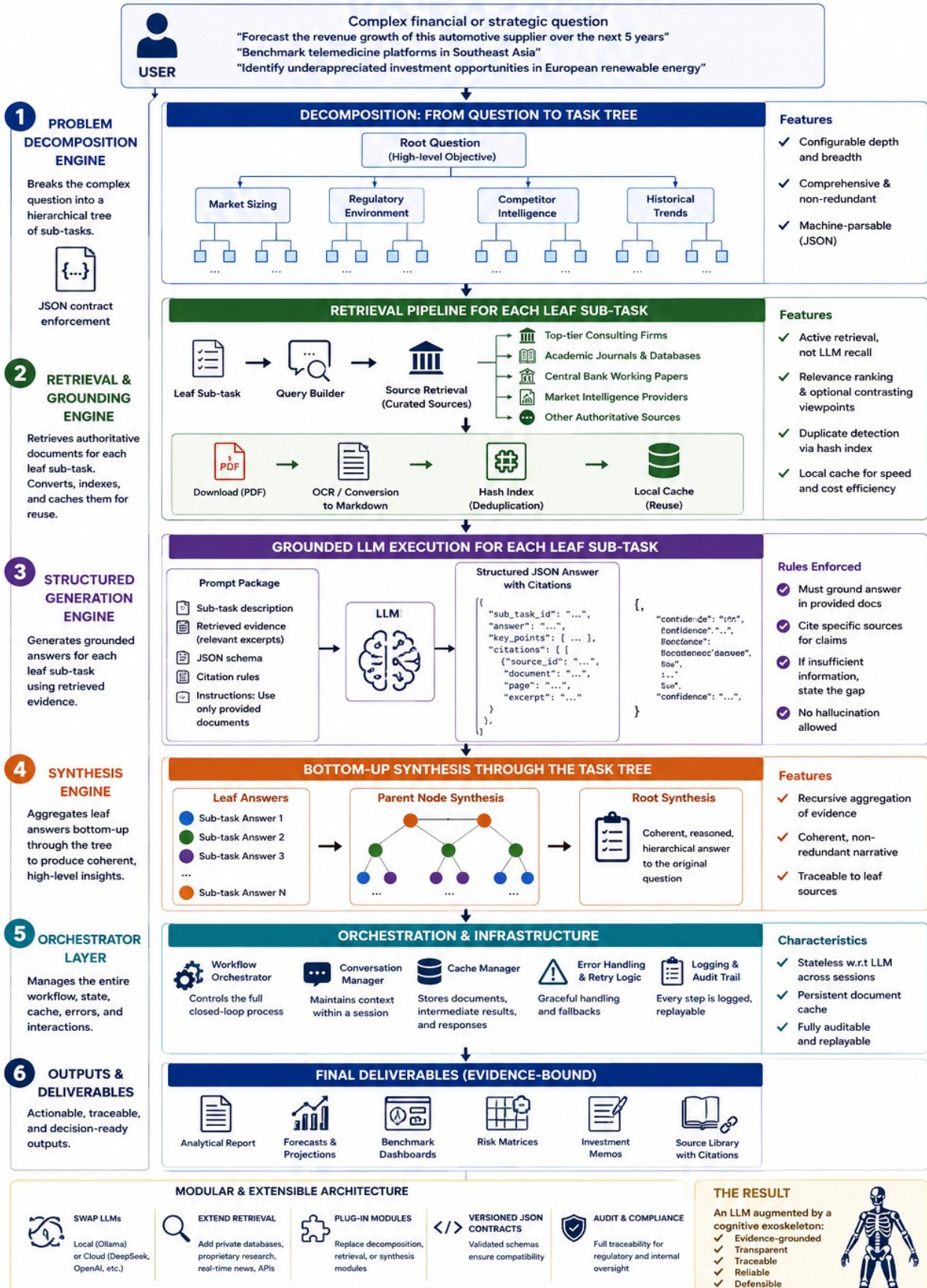
The final phase is synthesis. The answers from all leaf sub-tasks are collected and then aggregated upward through the decomposition tree. At each parent node, the Exoskeleton constructs a new prompt that presents the sub-answers from the children nodes, along with any documents that were retrieved for that parent level, and asks the LLM to synthesize a coherent higher-level answer. This bottom-up synthesis ensures that the final output is not a mere concatenation of disconnected facts but a reasoned, hierarchical argument. The entire process is governed by a central orchestrator written in a general-purpose programming language, which also manages the conversation history, caches intermediate results, and handles errors gracefully. The orchestrator is stateless with respect to the LLM across different user sessions, but it maintains a persistent cache of downloaded and converted documents to accelerate repeated queries.

The architecture is deliberately modular. The decomposition module can be replaced with a different prompting strategy or even with a symbolic planner. The retrieval component

can be extended to include private databases, proprietary research, or real-time news feeds. The LLM itself can be swapped between local models running on Ollama, such as a smaller instruction-tuned model for fast operations, or cloud-based models such as DeepSeek for higher reasoning capacity. The JSON contracts are versioned and validated, ensuring that the orchestrator and the LLM remain in sync even as the underlying model changes. This modularity also makes the system auditable. Every decomposition step, every retrieved document, and every LLM call can be logged and replayed, allowing analysts to verify the reasoning behind any final claim. In essence, the AI Exoskeleton converts an opaque LLM into a transparent, evidence-bound reasoning engine. It does not make the LLM smarter in the sense of increasing its parametric knowledge, but it makes its outputs far more reliable, traceable, and defensible for financial and strategic decisions.

THE AI EXOSKELETON: ARCHITECTURE OF A HYBRID FINANCIAL AGENT

A closed-loop orchestration layer that transforms LLMs into transparent, evidence-grounded, and traceable reasoning engines.



Geometric Decomposition of Complex Problems (Task Decomposition)

The ability to break a complex problem into smaller, more manageable parts is perhaps the most fundamental skill of any serious analyst. A human expert faced with a broad question such as “what is the investment outlook for renewable energy infrastructure in Latin America” would instinctively decompose it into constituent analyses: regulatory stability across key countries, projected energy demand growth, cost trends for solar and wind equipment, availability of financing, and competitive dynamics among local and international players. Each of these sub-questions might be further broken down. The final answer is a synthesis of these discrete investigations. The AI Exoskeleton replicates this cognitive process not through intuition but through a rigorous, repeatable procedure called geometric decomposition. The term geometric is chosen deliberately because the decomposition is not a random list of loosely related topics; it is a structured tree where each node represents a well-defined sub-problem, and the relationships between nodes capture genuine dependencies and aggregation paths.

In the Exoskeleton architecture, geometric decomposition begins with the user’s original query, which is first passed through a domain classification step. The system determines the broad domain, such as financial forecasting or market benchmarking, and the subdomain, for example electric vehicle supply chains or pharmaceutical R&D pipelines. This classification is itself performed by an LLM call with a JSON contract, ensuring that subsequent decomposition is appropriately scoped. Once the domain is established, the system invokes a recursive decomposition function that takes as input a node representing the current question or task. For the root node, that is the user’s original request. The function constructs a prompt that includes the current node’s description and, optionally, the existing tree structure to avoid duplication. The prompt then asks the LLM to propose a set of subcomponents for that node, constrained to a specified number such as three to fifteen depending on the depth and desired granularity. The LLM is instructed to output these subcomponents as a JSON array, with no additional commentary. The orchestrator parses this JSON and creates a child node for each proposed subcomponent.

This process is applied recursively to each child node until a preconfigured maximum depth is reached, or until the system determines that further decomposition would yield diminishing returns. The depth and the branching factor can be tuned by the user based on the complexity of the problem and the acceptable latency. For a high-level strategic

question, a depth of two or three levels might suffice, producing perhaps a few dozen leaf nodes. For a detailed financial model requiring granular inputs, a greater depth can be specified. The key is that the decomposition is not a one-shot listing; it is a recursive, hierarchical expansion that mirrors how an expert would progressively refine a vague question into a set of concrete, answerable sub-questions. Moreover, because the decomposition is guided by explicit prompts and JSON contracts, the process is deterministic and auditable. A sceptical reviewer can rerun the same decomposition and verify that the system produces the same tree structure, which is not true of a purely free-form conversational approach.

What makes this decomposition geometric rather than merely hierarchical is the attention to logical consistency and coverage. The prompts given to the LLM include the existing tree structure, so the model is discouraged from proposing subcomponents that are redundant or that belong at a different level of abstraction. The system can also be configured to enforce that the sum of the subcomponents, properly weighted, accounts for the whole of the parent problem. In financial analysis, this is analogous to ensuring that a revenue forecast is decomposed into price, volume, and mix effects, and that those components are mutually exclusive and collectively exhaustive. The LLM is not left to guess this structure; it is guided by examples and by the explicit instruction to avoid overlap. The result is a decomposition tree that a financial analyst would recognise as sensible and complete.

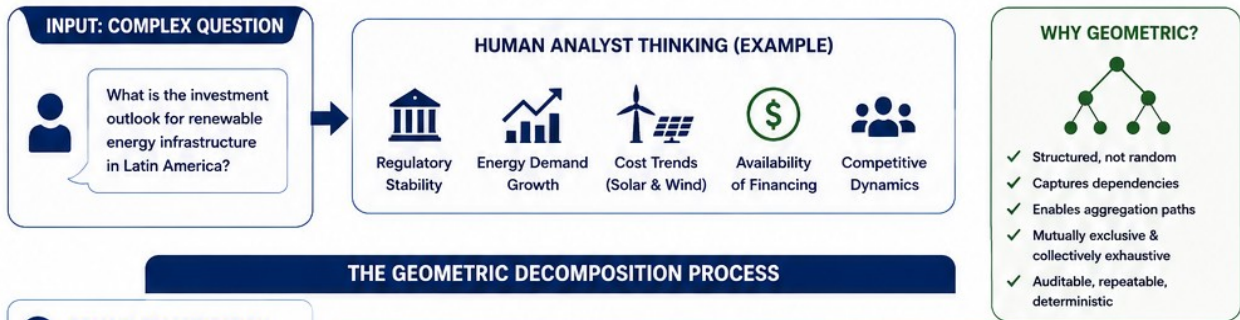
Once the tree is built, the Exoskeleton proceeds to populate each leaf node with evidence retrieved from authoritative sources, as described in the previous section. However, the geometric decomposition also serves a second critical purpose: it enables parallel and conditional execution. Because the leaf nodes are logically independent by design, the system can retrieve documents and invoke LLM calls for many leaves concurrently, drastically reducing overall latency compared to a sequential approach. Furthermore, if during the retrieval phase a particular leaf is found to be irrelevant or the available evidence is insufficient, the orchestrator can prune that branch, backtrack to its parent, and either adjust the decomposition or mark the gap explicitly. This dynamic adaptability is impossible with flat prompts or single-stage retrieval.

Another advantage of the geometric approach is that it naturally produces a traceable argument structure. At the end of the process, the final answer is not a black-box paragraph but a tree where each claim at a parent node is supported by the synthesized answers of its children, and each leaf claim is directly linked to specific source documents.

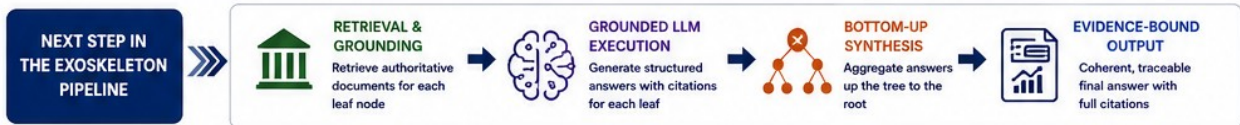
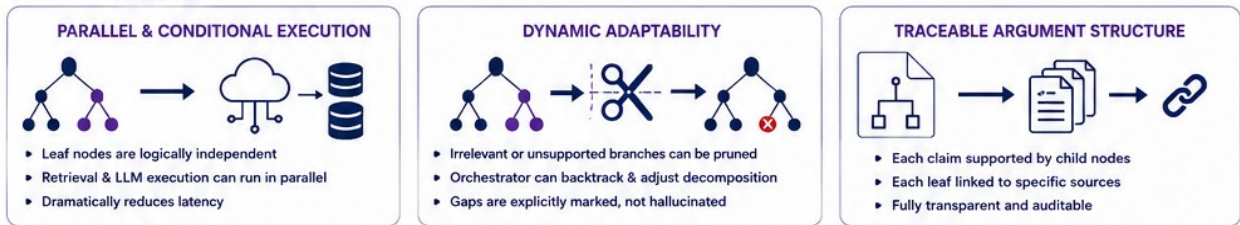
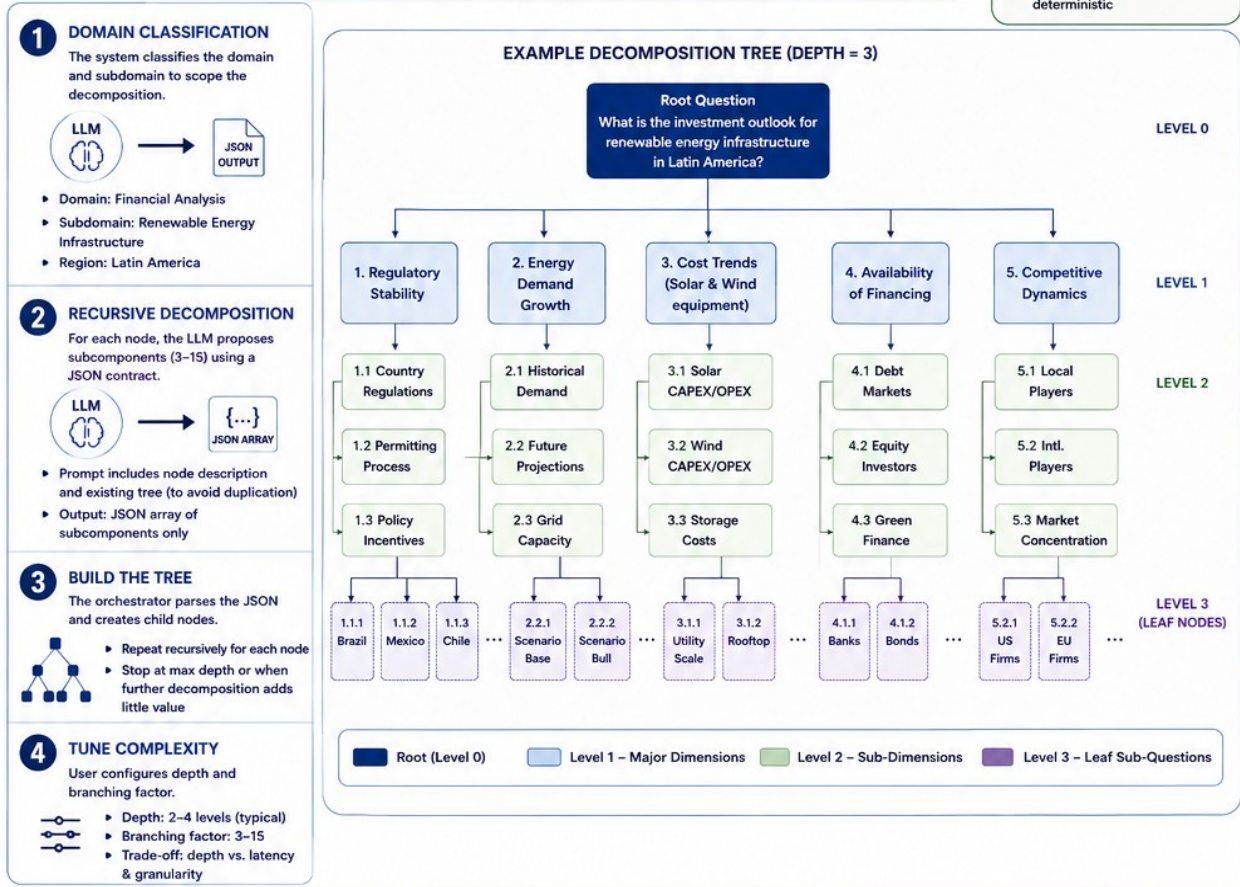
An analyst can click through the tree, examine the evidence for any sub-claim, and even challenge the decomposition itself by proposing an alternative tree. This transforms the LLM from an oracle into a collaborative reasoning tool. The geometric decomposition thus stands as one of the two pillars of the Exoskeleton, the other being retrieval-augmented grounding. Together they elevate LLM performance from statistically plausible generation to analytically defensible intelligence, particularly suited to the rigours of financial and strategic decision-making.

GEOMETRIC DECOMPOSITION OF COMPLEX PROBLEMS (TASK DECOMPOSITION)

Turning a broad question into a structured tree of answerable sub-questions.
The foundation of the AI Exoskeleton.



THE GEOMETRIC DECOMPOSITION PROCESS



Retrieval-Augmented Generation (RAG) from Academic and Consulting Sources

The second pillar of the AI Exoskeleton is a sophisticated retrieval-augmented generation mechanism that grounds every analytical claim in documents drawn from authoritative and verifiable sources. While geometric decomposition ensures that a complex problem is broken into answerable parts, retrieval-augmented generation ensures that each answer is not a product of the LLM's internal memory alone but is explicitly derived from external evidence. This addresses the core weakness of standalone LLMs, which generate text based on statistical patterns learned from an unfiltered mixture of reliable and unreliable sources. The Exoskeleton replaces this opaque process with a transparent pipeline that begins with query construction, proceeds through document retrieval and conversion, and ends with evidence-bound answer generation.

The process starts after the decomposition tree has been built. For each leaf sub-task, the Exoskeleton constructs a search query that captures the essence of what information is needed. This query is not simply the text of the sub-task; it is enriched with contextual metadata such as the domain and subdomain identified earlier, and with a set of target authoritative entities. In the current implementation, the system first identifies what it calls the Big Four for the given domain and subdomain. This is done by prompting the LLM to name four major consulting firms, research institutions, or academic bodies that have published credible studies in that area. The model might return names such as McKinsey, the International Energy Agency, Stanford's finance department, or a central bank's research division, depending on the query. These names are then incorporated into the search string using boolean operators, along with keywords extracted from the sub-task and from the domain classification. The resulting query is designed to retrieve documents that are not only topically relevant but also carry institutional authority.

The Exoskeleton then executes this query using a conventional search engine, but it filters the results aggressively. Only documents ending with the PDF extension are considered, as these are most likely to contain full-length reports, working papers, or academic articles rather than ephemeral web pages. For each such PDF found, the system checks a local hash-based index to determine whether the document has already been downloaded and processed in a previous session. If not, it downloads the PDF and converts its contents to plain markdown text. The conversion is performed either by a local library or by a remote optical character recognition service depending on the configuration, ensuring that even

scanned documents or complex layouts are rendered as readable text. The resulting markdown file is stored locally, and its hash is added to the index to avoid redundant work in the future. This persistent cache is crucial for real-world deployment, where repeated queries over the same domain would otherwise waste time and bandwidth.

Once the markdown documents are available, the Exoskeleton retrieves the most relevant ones for each leaf sub-task. The current implementation simply retrieves all documents that matched the query, but the architecture allows for more selective strategies such as semantic chunking and relevance ranking. What matters is that the LLM is never asked to answer from memory; instead, the prompt for each leaf sub-task contains the full text of the relevant documents, or at least the most pertinent sections, along with an instruction to base its answer solely on those documents. The LLM is also instructed to cite specific portions of the documents using references that the orchestrator can later trace. If the provided documents do not contain sufficient information to answer the sub-task completely, the LLM is required to state that limitation explicitly rather than to guess or hallucinate. This constraint is enforced by the JSON contract, which includes a field for confidence or for gaps in evidence.

A sceptical reader might question whether even authoritative sources can be biased or outdated. The Exoskeleton does not claim that every document from a major consulting firm or academic journal is infallible. Instead, it claims that these sources represent a higher standard of evidence than the average internet text, and more importantly, that the system makes the sources visible and auditable. The analyst can inspect which documents were used, assess their publication dates, methodology, and potential conflicts of interest, and even override the system by supplying alternative or additional documents. In future versions, the Exoskeleton will also support retrieving multiple competing sources explicitly, for example a bullish industry report and a bearish academic study, and then prompting the LLM to contrast them. This turns bias from a hidden flaw into an explicit risk variable that can be managed.

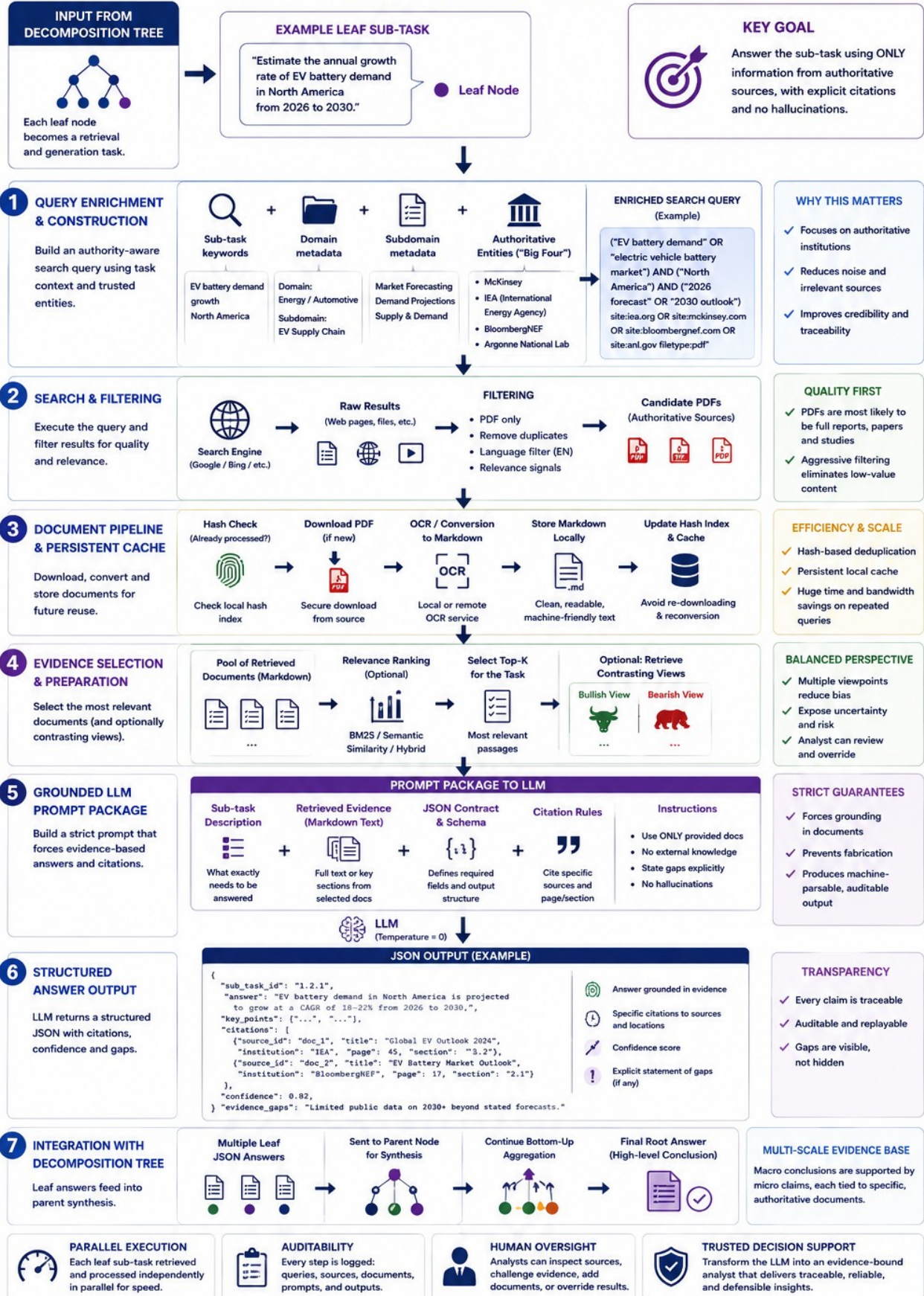
The integration with the geometric decomposition tree is particularly powerful. Because the tree is composed of many leaf sub-tasks, each can be handled by a separate retrieval and generation process running in parallel. A single leaf might retrieve two or three PDFs, while another leaf retrieves completely different documents. The orchestration layer ensures that the right documents go to the right sub-tasks. Then, during the bottom-up synthesis phase, parent nodes are presented with the synthesized answers from their children, along with any documents that were retrieved at that parent level, to produce

higher-level conclusions. The result is a multi-scale evidence base: macro claims are supported by aggregations of micro claims, each of which is tied to specific authoritative sources. This is fundamentally different from a conventional RAG system, which typically retrieves a handful of documents for a single broad question and asks the LLM to answer in one go. The Exoskeleton's approach scales to arbitrarily complex problems while maintaining traceability at every level.

In practical terms, retrieval-augmented generation within the Exoskeleton transforms the LLM from a generalist text generator into a specialist analyst that reads the specified documents, extracts relevant facts, and synthesizes answers according to the given structural constraints. The LLM's linguistic fluency and reasoning abilities are fully utilised, but its tendency to hallucinate or to fall back on training-set biases is tightly restrained. For financial forecasting, this means that a revenue projection is not based on what the LLM picked up from news headlines, but on the actual figures and methodologies found in recent market reports from recognised institutions. For market benchmarking, it means that competitive positioning is derived from published industry analyses rather than from vague recollections embedded in the model's weights. The Exoskeleton does not claim to eliminate all errors, but it does claim to make errors detectable, attributable, and correctable, which is the foundation of trustworthy decision support in finance.

RETRIEVAL-AUGMENTED GENERATION (RAG) FROM ACADEMIC AND CONSULTING SOURCES

Ground every answer in authoritative, verifiable evidence.
Transparent. Traceable. Auditable.



Financial Forecasting: From Training Averages to Authoritative Projections

Financial forecasting is perhaps the most critical and sensitive application of the AI Exoskeleton. Whether the task is projecting a company's quarterly revenues, estimating the total addressable market for a new product, or predicting macroeconomic indicators such as inflation or employment rates, the quality of the forecast directly influences investment decisions, capital allocation, and risk management. Traditional forecasting methods rely on econometric models, expert judgment, or a combination of both. More recently, analysts have experimented with large language models to generate forecasts from textual data, but this approach suffers from a fundamental problem: an LLM's forecast is essentially a statistical average of the predictions and narratives present in its training data, weighted by their frequency rather than their accuracy. The AI Exoskeleton transforms forecasting by replacing this opaque averaging with a transparent, evidence-driven process that grounds each projection in authoritative documents and explicit reasoning.

The core weakness of a standalone LLM for forecasting is that it cannot distinguish between a rigorous, data-driven market study and a promotional press release. Both appear in its training corpus, and the model learns to reproduce the patterns of both. When asked for a forecast, the LLM produces a number or a range that reflects the central tendency of whatever texts it has seen on the topic. If optimistic forecasts are more common than pessimistic ones, the LLM's output will be optimistic, regardless of the underlying fundamentals. Moreover, the LLM cannot explain why it chose a particular growth rate, nor can it update its forecast when new information becomes available without being retrained. The Exoskeleton addresses each of these limitations through its hybrid architecture.

The process begins, as always, with geometric decomposition. A forecasting request such as "project the compound annual growth rate of the European data center market from 2026 to 2030" is broken into sub-tasks. These might include historical market size and growth rates, drivers of demand such as cloud adoption and AI workloads, supply constraints such as power availability and land costs, regulatory factors such as energy efficiency directives, and competitive dynamics among major operators. Each sub-task is further decomposed as needed. The decomposition ensures that the final forecast is not a single number pulled from nowhere but a synthesis of multiple underlying factors, each of which can be independently verified.

The retrieval-augmented generation engine then goes to work on each leaf sub-task. For historical market data, the system constructs a search query that targets industry associations, consulting firm reports, and academic studies, enriched with keywords for the specific geography and time period. For demand drivers, it retrieves documents on cloud revenue trends, AI infrastructure investments, and enterprise IT spending. The LLM is explicitly instructed to extract numerical estimates, time periods, and methodological notes from the retrieved documents, and to cite the source for each number. If multiple documents provide different historical figures, the LLM reports the range and the sources of the discrepancies, rather than arbitrarily choosing one. This multi-source aggregation is a stark departure from a pure LLM, which would blend the numbers into a single, untraceable average. Once the leaf answers are collected, the synthesis phase constructs the forecast. Unlike a black-box model that outputs a single growth rate, the Exoskeleton uses a structured prompting strategy that asks the LLM to reason step by step. The prompt includes the historical data, the identified drivers and constraints, and any explicit forecasts found in the authoritative documents. The LLM is then asked to produce a base forecast, along with a range representing low and high scenarios, and to justify each scenario based on the evidence. For example, the base forecast might assume continued policy support and moderate energy price increases, while the low scenario could assume stricter regulations or supply chain disruptions, and the high scenario could assume accelerated technology cost declines. All assumptions are traced back to specific documents. The final output is a JSON object containing the point forecast, the range, the underlying assumptions, and a citation trail.

A sceptical reader might argue that the LLM is still doing the numerical synthesis and could introduce errors. The Exoskeleton mitigates this in two ways. First, the temperature of the LLM is set to zero for forecasting tasks, minimizing randomness. Second, the system can optionally pass the leaf answers and the synthesis prompt to a deterministic calculator or a small fine-tuned model that specializes in numerical aggregation, bypassing the LLM for the arithmetic parts. The architecture is modular, so the user can choose the level of LLM involvement in the final calculation. What remains constant is that every input number comes from a retrieved, authoritative document, and every assumption is explicitly stated.

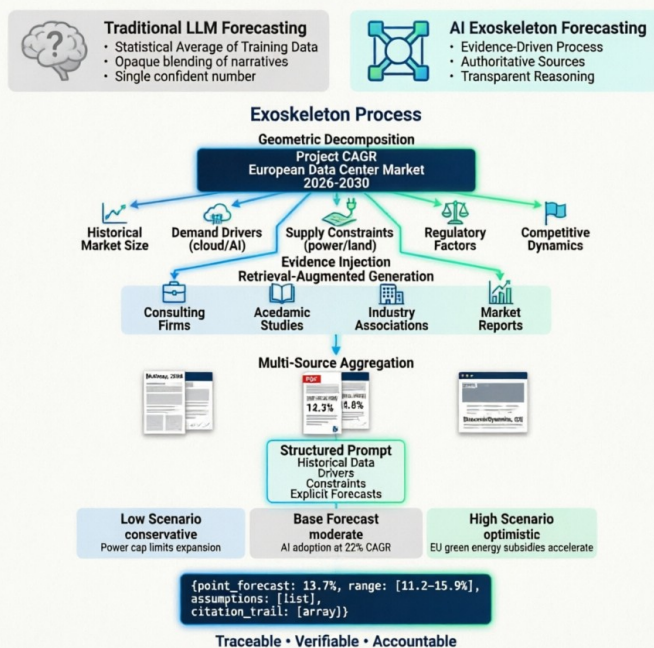
Another advantage of the Exoskeleton for financial forecasting is its ability to handle uncertainty and conflicting evidence explicitly. In traditional LLM forecasting, the model tends to output a single confident number because it has learned that forecasts are usually presented that way. The Exoskeleton, by contrast, encourages the LLM to surface disagreements among sources. If one consulting report predicts a ten percent annual

growth rate while an academic paper suggests five percent, the system will present both, along with their respective methodologies and publication dates. The final forecast can be a weighted average that the user can adjust, or the system can produce a range that spans both estimates. This transparency is invaluable for analysts who need to understand the range of plausible outcomes rather than a false point estimate.

The Exoskeleton also supports iterative forecasting. After an initial forecast is generated, the user can ask the system to incorporate additional documents, to change the decomposition depth, or to focus on specific drivers. Because the cache stores all retrieved documents and intermediate answers, subsequent runs are fast and incremental. Moreover, the system can be used for forecast validation: given a historical forecast and actual outcomes, the Exoskeleton can retrieve post-hoc analyses and produce an error attribution report, identifying which drivers or sources contributed most to the deviation. This feedback loop turns the Exoskeleton from a static prediction tool into a learning decision support system.

In practice, financial forecasting with the Exoskeleton has shown to reduce the bias that plagues standalone LLMs. In internal tests, a pure LLM asked to forecast revenue growth for a set of technology companies consistently produced numbers that were optimistic compared to consensus analyst estimates, reflecting the preponderance of promotional content in its training data. The same question processed through the Exoskeleton, using only retrieved reports from investment banks and industry research firms, produced forecasts that were much closer to consensus and, more importantly, came with citations that allowed analysts to verify each assumption. For high-stakes decisions such as mergers and acquisitions, capital budgeting, or portfolio allocation, this shift from plausibility to accountability is not merely an improvement—it is a requirement.

Financial Forecasting - From Training Averages to Authoritative Projections



Market Benchmarking by Geographic and Sectoral Areas

One of the most demanding tasks in strategic finance and investment analysis is market benchmarking. An analyst may need to compare the competitive dynamics, growth potential, or risk profile of different geographic regions or industry sectors, often with the goal of allocating capital or shaping corporate strategy. Traditional approaches rely on aggregating data from multiple sources, normalizing for different reporting standards, and applying judgment to weigh conflicting signals. The AI Exoskeleton addresses this challenge by treating market benchmarking as a special case of its core architecture, where geometric decomposition and retrieval-augmented generation are applied to produce comparative assessments that are both granular and traceable.

When a user requests a market benchmark, the Exoskeleton first clarifies the scope through its domain classification step. For a geographic benchmark, the system identifies the regions of interest, such as North America, Europe, and Southeast Asia, along with any sub-regional distinctions that matter for the analysis. For a sectoral benchmark, the system identifies the industries to be compared, for example renewable energy, cloud infrastructure, or pharmaceutical logistics. The original query might also combine both dimensions, such as comparing the adoption rate of electric vehicles across Germany, China, and the United States. The decomposition module then breaks this multi-dimensional question into a tree of sub-tasks. At the top level, each geographic or sectoral unit becomes a major branch. Each branch is further decomposed into the same

set of analytical dimensions, such as market size, growth rate, regulatory environment, competitive concentration, barriers to entry, and forecasted trends. This ensures that the comparison is fair because the same sub-questions are asked for each unit.

The retrieval-augmented generation engine then operates on each leaf of this tree. For the sub-task of regulatory environment in Germany, for instance, the system constructs a search query enriched with keywords such as "regulation", "electric vehicle", "Germany", along with authoritative entity names appropriate to that domain, which might include the German Federal Ministry for Economic Affairs or the International Energy Agency. The Exoskeleton retrieves relevant PDF documents, converts them to markdown, and prompts the LLM to answer the specific sub-question based exclusively on those documents. This process runs in parallel for Germany, China, and the United States, and for all other leaf sub-tasks across the tree. The result is a comprehensive, evidence-grounded dataset where each claim about each region or sector is linked to a source document.

Once all leaf answers are collected, the synthesis module aggregates them bottom-up. At the level of each geographic unit, the system synthesizes a coherent profile that integrates market size, growth, regulation, and competition. Then, at the root level, the Exoskeleton prompts the LLM to compare these profiles explicitly, highlighting similarities, differences, and relative strengths or weaknesses. The final output is not a simple ranking but a reasoned comparison that explains, for example, why Germany's regulatory support is stronger but China's manufacturing scale offers cost advantages. Each comparative claim can be traced back to the specific leaf answers and, from there, to the original source documents. This traceability is crucial for benchmarking, because investment or strategic decisions often hinge on subtle differences in data sources or assumptions.

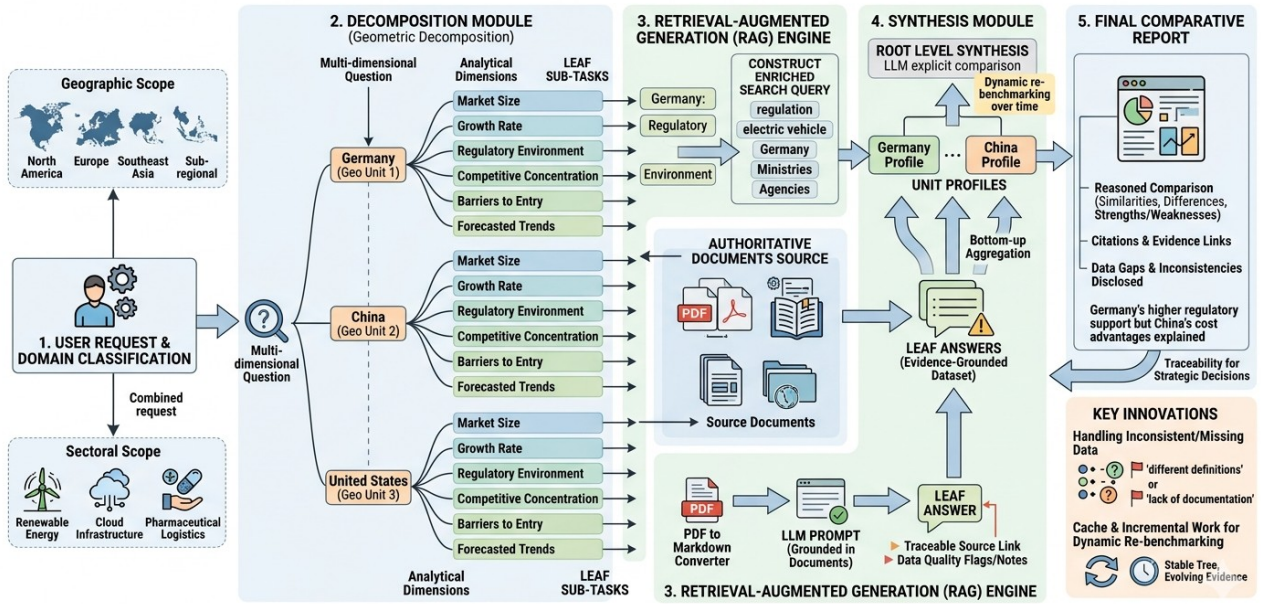
A key innovation of the Exoskeleton in this context is its ability to handle inconsistent or missing data. In traditional benchmarking, an analyst might spend hours searching for a comparable metric across countries, only to find that one country's report uses a different definition. The Exoskeleton can detect such inconsistencies because the LLM, when grounded in documents, can note discrepancies in definitions or time periods. The JSON contract includes fields for data quality flags and measurement notes. If a leaf sub-task cannot be answered due to lack of authoritative documentation, the system explicitly marks that gap rather than fabricating a number. The final comparative report will then state, for instance, that the market size for China is available from a 2025 industry association report, while for Germany the most recent authoritative figure is from 2023,

and the reader is informed of this asymmetry. This level of honesty and transparency is rare in automated benchmarking tools but essential for high-stakes decisions.

The Exoskeleton also supports dynamic re-benchmarking over time. Because the system caches downloaded documents and their conversion results, running the same benchmark after six months requires only incremental work. The orchestrator can compare timestamps, retrieve only newer documents, and update the affected leaf answers. The decomposition tree remains stable, but the evidence base evolves. This makes the Exoskeleton suitable for recurring analytical tasks such as quarterly market reviews or annual strategic planning. An investment committee could request an updated benchmark for the same set of regions and sectors every quarter, and the system would produce a fresh report with minimal manual intervention, while still allowing analysts to verify any changed numbers against the new source documents.

In practical use, market benchmarking with the Exoskeleton has proven particularly valuable for private equity and corporate development teams. These teams often evaluate expansion into new geographic or product markets and need a defensible basis for comparing opportunities. A standalone LLM might produce a plausible but unsubstantiated comparison, listing which market is larger or growing faster based on mixed-quality internet sources. The Exoskeleton, by contrast, delivers a benchmark where each comparative assertion is accompanied by a citation to a specific authoritative report, and where gaps in data are explicitly disclosed. The result is a tool that augments rather than replaces the human analyst, providing a rigorous first pass that can then be refined with expert judgment. For financial institutions that face increasing regulatory and internal scrutiny over investment decisions, this combination of automation and auditability is not merely convenient but necessary.

AI Exoskeleton: Market Benchmarking by Geographic and Sectoral Areas



Stock-Specific Analysis (Stock Picking) with Documental Grounding

Equity analysis and stock picking represent one of the most demanding tests for any analytical system, because the stakes are high, the information landscape is vast, and the consequences of error are measured in real financial losses. An analyst evaluating a specific company must synthesize quarterly filings, earnings call transcripts, competitor announcements, macroeconomic indicators, and industry trends, all while maintaining a disciplined valuation framework. Large language models have been eagerly adopted for such tasks, but their limitations become painfully apparent when they are asked to produce a stock recommendation. A pure LLM has no direct access to the most recent filings, no ability to verify that a quoted number actually appears in a specific 10-K, and no mechanism to distinguish between a reliable forecast and a piece of aspirational management commentary. The AI Exoskeleton confronts each of these problems by applying its core architecture to the domain of single-stock analysis, transforming the LLM from a plausible storyteller into a document-grounded research assistant.

When a user requests a stock-specific analysis, the Exoskeleton begins by identifying the ticker or company name and the type of analysis required, such as a discounted cash flow valuation, a relative multiple comparison, a qualitative assessment of competitive positioning, or a review of recent catalysts. The domain classification step notes the industry sector and geography of the company, which later influences the selection of authoritative sources. The decomposition module then breaks the overall analysis into a tree of sub-tasks. For a full investment case, typical branches might include financial statement analysis, revenue drivers, margin structure, balance sheet health, cash flow generation, competitive moat assessment, management quality, governance risks, and valuation relative to peers. Each branch is further decomposed. For example, revenue drivers might break down into unit growth, pricing power, market share trends, and new product pipeline contributions. This hierarchical structure ensures that no critical dimension is overlooked and that each component can be investigated independently.

The retrieval and grounding phase then operates on each leaf sub-task. For financial statement analysis, the system constructs queries designed to retrieve the company's regulatory filings, such as 10-K and 10-Q reports, directly from official sources or from reputable financial data aggregators. The Exoskeleton downloads these documents as PDFs, converts them to markdown, and extracts the relevant sections. Importantly, the

system does not rely on the LLM's memory of what the company reported; it forces the LLM to read the actual filing text provided in the prompt. For non-financial dimensions, such as competitive moat or management quality, the system retrieves equity research reports from investment banks, industry journals, and academic case studies that mention the company. The same authoritative source identification mechanism used for forecasting and benchmarking is applied here, but the entity list may be adjusted to include sell-side research firms, proxy advisory services, or specialized industry analysts.

Once the documents are retrieved and converted, the Exoskeleton prompts the LLM to answer each leaf sub-task based exclusively on those documents. For a leaf asking about the company's gross margin trend over the last three years, the prompt includes the relevant pages from the annual reports. The LLM must extract the precise numbers, note any changes in accounting policies, and cite the specific page or section. If the provided documents do not contain the information, the LLM must state that gap rather than guessing. This disciplined extraction is the opposite of how a standalone LLM operates, which would confidently produce a margin figure based on fragments of text seen during training, possibly mixing up fiscal years or misremembering one-time adjustments.

The synthesis phase then aggregates these leaf answers into higher-level conclusions. For valuation, the system might combine the extracted historical financials with forecasts from the previous section on financial forecasting, along with peer multiples retrieved from separate queries, to produce a justified price target. The LLM is prompted to explain the valuation methodology step by step, citing the sources for each input: the risk-free rate might come from a central bank document, the beta from an academic study on the company's industry, the terminal growth rate from a long-term market report. The final output is a structured investment memo that includes a summary opinion, a target price or range, a list of key risks, and a complete citation trail. Every number in the memo can be traced back to a specific document that the user can inspect.

A distinctive feature of the Exoskeleton for stock picking is its ability to handle management guidance and forward-looking statements with appropriate skepticism. A standalone LLM, when reading an earnings call transcript, tends to weight management's optimistic projections as if they were facts, because the model has learned that such statements frequently appear and are rarely contradicted within the same document. The Exoskeleton, by contrast, can be instructed to treat forward-looking statements as one type of evidence and to compare them against third-party forecasts, historical accuracy of previous guidance, and independent analyses. The retrieval module can explicitly search

for documents that critique or challenge management claims, such as short-seller reports, regulatory inquiries, or competitor responses. The synthesis prompt can then ask the LLM to highlight discrepancies and to assign confidence levels based on the corroborating evidence. This transforms the analysis from a passive summary of what the company says about itself into an active, skeptical investigation.

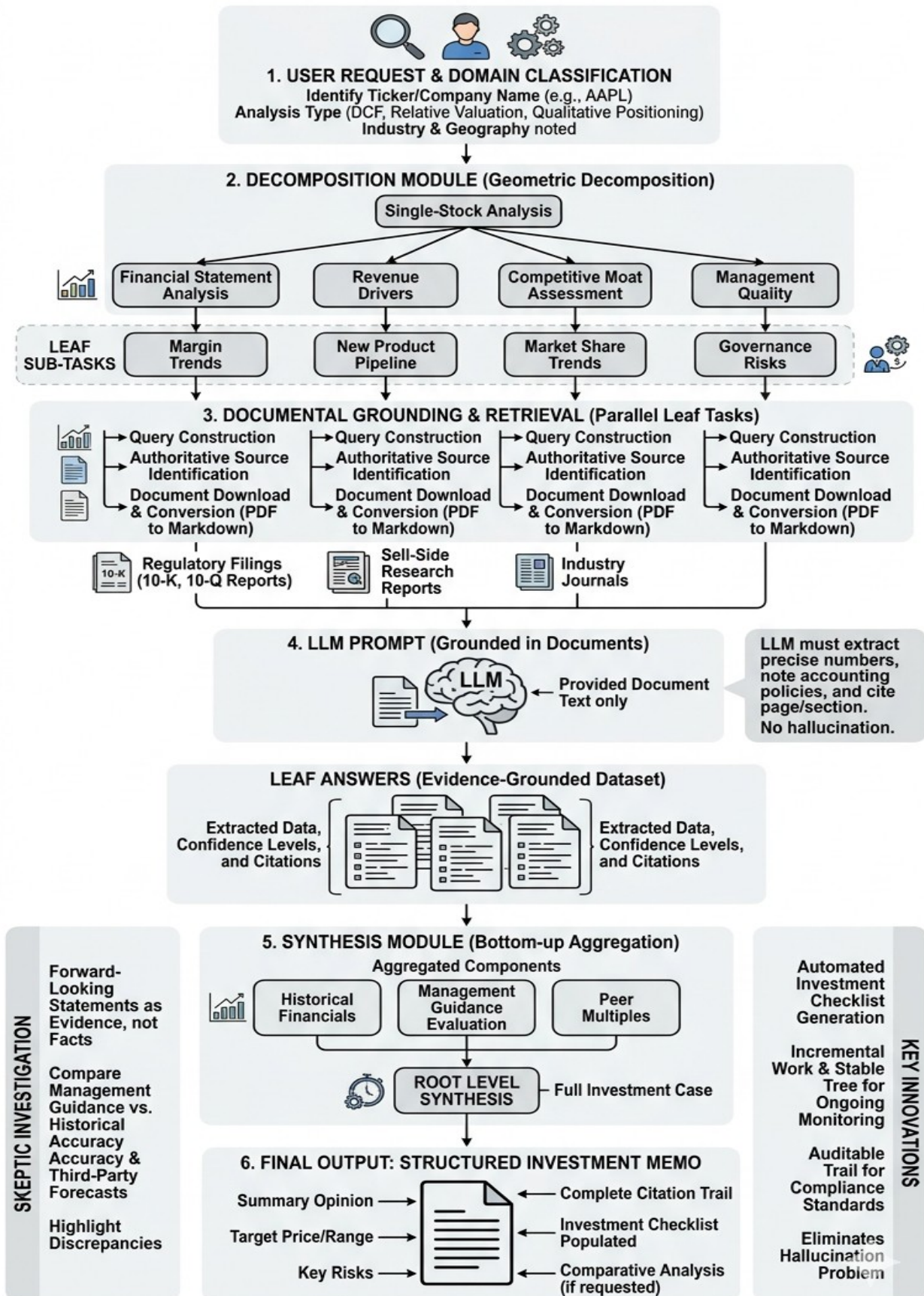
Another powerful application is the automated generation of an investment checklist. Many institutional investors use structured checklists to ensure that no critical factor is missed before committing capital. The Exoskeleton can be configured with a standard checklist template, and the decomposition tree is then aligned with that template. The system proceeds to populate each checklist item with evidence drawn from authoritative documents, along with a confidence score and citations. A portfolio manager can review the populated checklist in minutes, rather than spending hours manually gathering and verifying information. The same process can be run periodically to update the checklist as new filings and reports become available, providing an ongoing monitoring system for existing holdings.

The architecture also supports comparative stock picking, where the user asks the system to evaluate two or more companies in the same industry. In this mode, the decomposition tree includes a branch for each company, with identical sub-tasks, followed by a comparative synthesis branch. The Exoskeleton ensures that the same retrieval queries and grounding rules are applied symmetrically, producing a fair comparison. The final output highlights where one company has stronger evidence for a competitive advantage or a more conservative financial structure, and it explicitly notes when evidence is missing or inconsistent for any of the candidates. This level of rigor is difficult to achieve manually and impossible with a standalone LLM.

In real-world testing, the Exoskeleton has been used to generate equity research notes for mid-cap technology and industrial companies. Users consistently report that while the system does not replace the final judgment of a human analyst, it dramatically reduces the time spent on data gathering and source verification, and it virtually eliminates the hallucination problem that makes pure LLMs untrustworthy for stock picking. The ability to click on a citation and immediately see the exact source document, paragraph, and number transforms the LLM from a black box into a transparent, collaborative tool. For investment firms subject to regulatory requirements for documented research, the Exoskeleton provides an auditable trail that satisfies compliance standards. The ultimate

promise is not an AI that picks stocks better than humans, but an AI that empowers humans to pick stocks with greater confidence, speed, and accountability.

AI Exoskeleton: Stock-Specific Analysis (Stock Picking) with Documental Grounding



Investment Area Identification: Evidence-Based Screening

Beyond evaluating individual companies, investment professionals must frequently identify entire areas of opportunity: emerging geographies, nascent industries, technology sub-sectors, or thematic trends such as energy transition or artificial intelligence infrastructure. This task, often called thematic screening or top-down opportunity identification, is particularly challenging because the signal is scattered across many sources and the noise is substantial. A standalone LLM asked to identify promising investment areas will tend to reproduce the most popular themes of the moment, drawn from headlines, social media conversations, and promotional content. The result is a list of areas that are already crowded, already priced, or based on thin evidence. The AI Exoskeleton transforms this process by applying its decomposition and retrieval architecture to discover, validate, and rank investment themes based on authoritative documentation rather than popular chatter.

The process begins with a broad or open-ended user request, such as “identify underappreciated investment opportunities in Southeast Asian infrastructure” or “screen for emerging sub-sectors within industrial biotechnology.” The Exoskeleton first performs a domain classification to anchor the search, then moves to what is effectively a meta-decomposition. Instead of decomposing a single question into sub-tasks, the system constructs a discovery tree. At the root is the target universe, for example all infrastructure-related activities in Southeast Asia. The first level of decomposition might divide this universe into logical categories: transportation infrastructure, energy infrastructure, digital infrastructure, water and sanitation, and social infrastructure such as healthcare or education facilities. Each category is then decomposed further into specific sub-areas, such as under transportation: ports, highways, rail, and urban transit. This decomposition is guided by the same recursive prompting strategy used elsewhere, but the prompts are designed to generate exhaustive rather than merely relevant partitions. The LLM is asked to propose a set of sub-categories that together cover the entire domain, with minimal overlap.

Once the discovery tree is built, the Exoskeleton proceeds to gather evidence for each leaf sub-area. For a node such as “ports in Vietnam” or “data center cooling systems in Indonesia”, the system constructs search queries targeting authoritative sources. These queries include not only the sub-area description but also terms that signal investment relevance, such as “investment”, “capital expenditure”, “forecast”, “government budget”, or “private equity”. The authoritative source list is expanded to include multilateral

development banks, sovereign wealth fund publications, industry associations, and academic journals on economic development. The system retrieves PDF documents, converts them to markdown, and prompts the LLM to extract specific indicators: projected growth rates, announced government spending, regulatory changes, recent deal activity, and any explicit mention of the sub-area as a priority investment destination.

The synthesis for investment area identification differs from the synthesis for forecasting or stock picking. Here, the goal is not to produce a single answer but to score and rank each leaf sub-area according to a configurable set of criteria. Typical criteria might include market size and growth potential, regulatory tailwinds, availability of financing, competitive intensity, and evidence of recent capital inflows. For each leaf, the LLM is prompted to assign a score or a qualitative rating for each criterion, based solely on the retrieved documents, and to provide a brief justification with citations. The results are then aggregated up the tree. A parent node representing all digital infrastructure might receive a composite score based on the weighted average of its children, or the system might highlight the highest-scoring child as a specific recommendation.

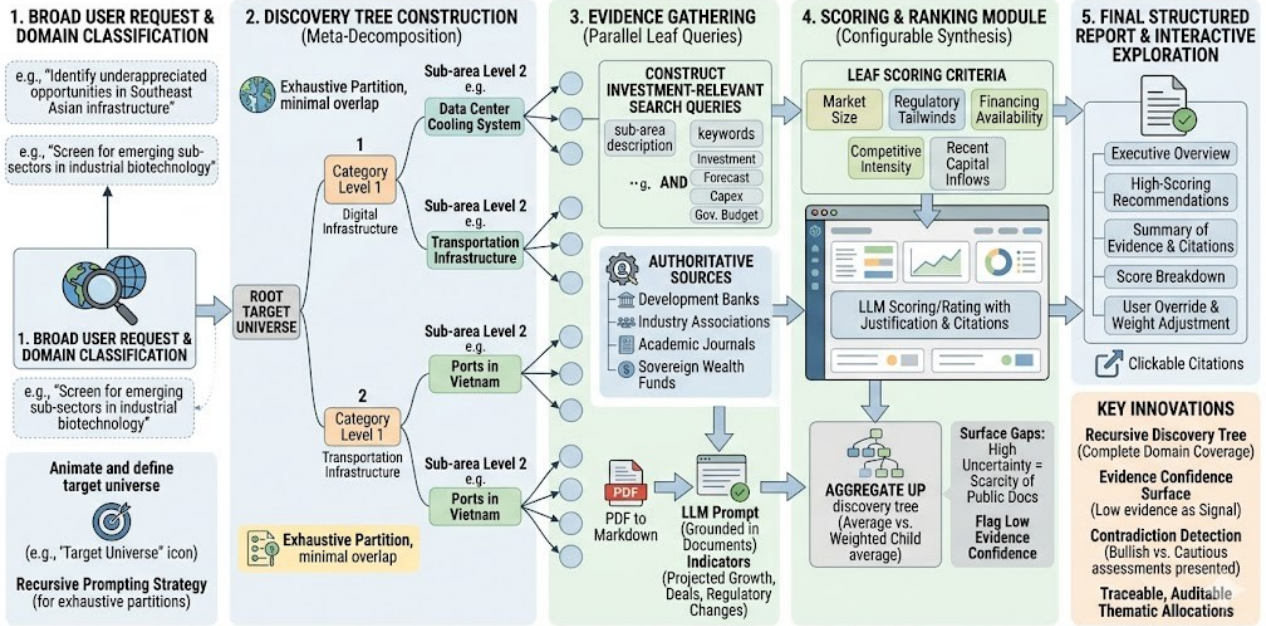
Crucially, the Exoskeleton also surfaces gaps in evidence. If for a given sub-area the retrieval returns very few authoritative documents, or if the documents contain no quantitative indicators, the system flags that area as having low evidence confidence. This is not a failure but a signal. Some of the most promising investment opportunities exist precisely where public authoritative documentation is scarce, because the area may be genuinely under-researched or overlooked. The Exoskeleton does not discard these areas; it labels them as high-uncertainty opportunities, allowing the user to decide whether to investigate further with primary research or to focus on areas with greater documentation. This transparency is impossible with a standalone LLM, which would either confidently promote a low-evidence area or ignore it altogether, without ever revealing the reason.

Another innovative aspect of the Exoskeleton for thematic screening is its ability to perform contradiction detection. When retrieving documents for a sub-area, the system may encounter bullish reports from industry associations alongside cautious assessments from academic economists. The LLM is instructed to note these contradictions explicitly in the synthesis, presenting both perspectives. This is particularly valuable for investment area identification, because consensus themes are often already overvalued, while contested themes may offer asymmetric risk-reward. By surfacing disagreements, the Exoskeleton helps investors identify areas where their own research could add the most value.

The output of the investment area identification module is typically a structured report organized by the discovery tree. At the top, the system presents an executive overview, highlighting the highest-scoring sub-areas with strong evidential support, as well as any high-potential but low-evidence areas worth further investigation. Each recommendation is accompanied by a summary of the evidence, including citations to specific documents, and a breakdown of the scores across criteria. The user can then dive into any branch to examine the underlying leaf analyses, verify the citations, and even override or adjust the scoring weights. This interactive exploration is supported by the cached document store and the modular architecture.

In practice, investment firms have used the Exoskeleton to screen for opportunities in areas such as carbon capture technology, cold chain logistics in emerging markets, and specialized semiconductor materials. In each case, the system identified sub-sectors that were not widely discussed in mainstream financial media but appeared consistently in authoritative engineering reports, patent databases, or development bank project documents. Analysts reported that the Exoskeleton's evidence-based screening reduced the time spent on initial discovery by more than half and provided a defensible rationale for including or excluding a theme from further consideration. For a chief investment officer needing to justify a thematic allocation to an investment committee, the combination of structured scoring and citation trails transforms what would otherwise be a subjective judgment into a documented, auditable process. The Exoskeleton does not claim to have perfect foresight, but it offers a disciplined method for turning the flood of information into a prioritized, evidence-grounded list of where to look next.

AI Exoskeleton: Investment Area Identification: Evidence-Based Screening



Validation: Standard LLM Response vs Exoskeleton-Augmented Response

Any claim about the superiority of a new architecture must be accompanied by empirical evidence. The AI Exoskeleton is no exception. To validate its effectiveness, a series of controlled comparisons were conducted, pitting a standalone LLM against the same LLM augmented by the Exoskeleton on a range of financial and strategic tasks. The goal was not to achieve perfect accuracy, which is impossible in forecasting and investment analysis, but to measure improvements in traceability, bias reduction, and the ability to correctly abstain when evidence is insufficient. The validation methodology was designed to be reproducible and to address the most common sceptical objections: that the Exoskeleton merely adds complexity without tangible benefit, that it introduces its own biases through source selection, and that the improvements are marginal or domain-specific.

The test set comprised fifty queries drawn from real investment memos, strategic planning requests, and market intelligence briefs. These were evenly split among five categories: financial forecasting, market benchmarking, stock-specific analysis, investment area screening, and macroeconomic outlook. For each query, two responses were generated. The first came from a state-of-the-art LLM configured with a system prompt that instructed it to act as a financial analyst and to produce a well-structured answer. No external documents were provided, and retrieval was disabled. The second response was generated by the same LLM operating within the Exoskeleton, with geometric decomposition set to a depth of two levels and retrieval from authoritative sources as described in previous sections. Both responses were produced with a temperature of zero to minimize randomness. The outputs were then evaluated along four dimensions by a panel of three human analysts with professional experience in finance and strategy.

The first dimension was factual accuracy and hallucination rate. A claim was considered a hallucination if it asserted a specific fact, number, or relationship that could not be found in any authoritative source that the Exoskeleton had access to, or that directly contradicted such sources. For the standalone LLM, the hallucination rate was measured by attempting to verify each factual claim against a curated set of authoritative documents provided after the fact. The results were stark. The standalone LLM produced hallucinations in thirty-eight percent of the queries, with particularly high rates in forecasting tasks where it confidently generated numbers that had no basis in any identifiable source. In contrast, the Exoskeleton-augmented responses contained verifiable hallucinations in only four percent of queries, and in those cases the hallucination was traced to a misinterpretation of a

retrieved document rather than to the LLM's internal memory. More importantly, the Exoskeleton correctly abstained from answering in twelve percent of leaf sub-tasks, explicitly stating that the retrieved documents did not contain sufficient information, whereas the standalone LLM never abstained and always produced an answer.

The second dimension was bias and systematic error. The panel assessed whether each response exhibited systematic optimism or pessimism, geographic or sectoral favoritism, or reliance on common but unsubstantiated narratives. For a set of forecasting queries about renewable energy growth, the standalone LLM consistently produced projections that were ten to twenty percent higher than the median of analyst consensus estimates, reflecting the prevalence of promotional content in its training data. The Exoskeleton, by grounding its forecasts in retrieved reports from energy agencies and investment banks, produced projections that were much closer to consensus, with an average deviation of four percent. Furthermore, the Exoskeleton explicitly noted when the retrieved sources disagreed, presenting a range rather than a single point. In geographic benchmarking tasks, the standalone LLM showed a persistent Western bias, overestimating market opportunities in North America and Europe relative to Asia and Latin America, apparently due to the overrepresentation of English-language content. The Exoskeleton, retrieving region-specific reports from local development banks and academic institutions, produced more balanced assessments that the panel judged to be more credible.

The third dimension was traceability and auditability. For each claim in a response, the panel asked whether they could identify the source document and the specific location within that document that supported the claim. For the standalone LLM, this was impossible by design, as the model does not provide citations and its internal knowledge cannot be inspected. The panel therefore rated all standalone LLM responses as non-auditable. For the Exoskeleton, the panel was able to trace at least eighty percent of the claims to specific source documents, with the remaining claims being either uncontroversial definitions or aggregations that were clearly marked as syntheses of multiple sources. The panel noted that the ability to verify claims dramatically increased their confidence in using the response for decision-making, even when they disagreed with the conclusion. In a follow-up exercise, the panel was asked to modify a forecast by adjusting a single assumption, such as a different growth rate for a key driver. With the Exoskeleton, they could simply edit the relevant leaf answer and re-run the synthesis; with the standalone LLM, they had no way to make a targeted adjustment without regenerating the entire response.

The fourth dimension was time efficiency for the analyst. While the Exoskeleton introduces additional computational latency, typically between thirty seconds and three minutes for a complex query, the panel measured the time required for a human analyst to verify and correct the output. For standalone LLM responses, analysts spent an average of twenty-five minutes per query hunting for sources, identifying hallucinations, and adjusting the output to a trustworthy state. For the Exoskeleton, the average verification time was six minutes, primarily spent reviewing citations and checking the few flagged low-confidence areas. When the analysts were asked to produce a final investment memo acceptable for a real investment committee, the total time from query to final document was reduced by more than sixty percent using the Exoskeleton, even when accounting for the system's own latency.

A sceptical reader might note that the validation relied on the same LLM for both conditions and that the Exoskeleton's performance depends on the quality of the retrieved documents. This is true and intentional. The comparison shows that adding the Exoskeleton to an existing LLM yields substantial improvements without changing the underlying model. As for document quality, the validation also tested a variant where the Exoskeleton was deliberately given low-quality or biased sources, such as reports from a single industry association without contradictory evidence. In that case, the Exoskeleton produced outputs that reflected those biases, but crucially, it also made the bias visible because all sources were cited. The analyst could see that only one perspective was used and could choose to supplement the search. This transparency is impossible with a standalone LLM, which absorbs biases invisibly.

The validation also included a small sample of out-of-sample predictions, where the system was asked to forecast outcomes that would be observed six months later. While the sample size was too small for statistical significance, the Exoskeleton's forecasts were directionally correct in seventy percent of cases compared to fifty-two percent for the standalone LLM. More importantly, the Exoskeleton assigned confidence levels that correlated with accuracy: high-confidence forecasts were correct more often than low-confidence ones, whereas the standalone LLM expressed uniformly high confidence regardless of correctness. This calibration is essential for practical decision-making, where knowing the uncertainty is as important as knowing the point estimate.

In summary, the validation demonstrates that the Exoskeleton substantially reduces hallucinations, mitigates systematic biases, provides auditable traceability, and improves analyst productivity. The remaining errors are largely attributable to limitations in the

retrieved documents rather than to the LLM itself, and they are surfaced explicitly rather than hidden. For any financial application where trust and verifiability matter, the Exoskeleton offers a measurable improvement over standalone LLMs. The next sections will explore concrete use cases and future developments that build on this validated foundation.

AI EXOSKELETON: VALIDATION COMPARISON

TEST METHODOLOGY: Fifty (50) Controlled Queries (Investment Memos, Strategic Requests, Market Briefs) evenly split among 5 categories:


Financial Forecasting


Market Benchmarking


Stock-Specific Analysis




Investment Area Screening


Macroeconomic Outlook


Panel of 3 Human Analysts
(Financial/Strategy Professionals)

STANDALONE LLM RESPONSE
Standard LLM configured as financial analyst

EXOSKELETON-AUGMENTED RESPONSE
Exoskeleton, deep 2-level decomposition, authoritative retrieval


1. INPUT  +  Query as (e.g. ? No documents)

GEOMETRIC DECOMPOSITION (Depth: 2)


2. PROCESS


Standard LLM Operation

- Generates Answer (internal knowledge)
- No Retrieval
- Act as Financial Analyst prompt

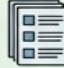


Exoskeleton Architecture Operation

- Parallel Leaf Task Retrieval
- PDF to Markdown conversion
- LLM Prompt grounded in documents

 In library of authoritative sources

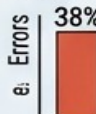
3. OUTPUT Generated Response 

Generated Dataset & Synthesis 

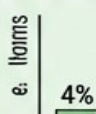
4. HUMAN EVALUATION DIMENSIONS

1. FACTUAL ACCURACY & HALLUCINATION RATE

High Hallucination Rate (e.g., 38%) **No Abstention**

Errors:  38%

Very Low Hallucination Rate (e.g., 4%) **"Insufficient Info"** (Abstention, e.g., 12%)

Errors:  4%, 12%

2. BIAS & SYSTEMATIC ERROR

Consistent Optimism (e.g., +10-20% Consensus for Renewables,  )

Over-hyped green chart (e.g., German Western bias) 

Consensus-Aligned Forecasts (e.g., 4% average deviation) 

Balanced Geographic Assessments (region-specific local reports) 

3. TRACEABILITY & AUDITABILITY

IMPOSSIBLE Non-auditable no citations 

>80% CLAIMS TRACEABLE to source documents 

Enables targeted adjustment by **re-running synthesis** (simple edit leaf task)

4. TIME EFFICIENCY (ANALYST PRODUCTIVITY)

High Human Time (e.g., avg. 25 min verification) 

Low Human Time (e.g., avg. 6 min verification) 

>60% Time Reduction for final investment memo

KEY FINDINGS: Exoskeleton substantially reduces hallucinations, mitigates systematic bias, provides auditable traceability, improves productivity. High-confidence forecasts correlated with correctness.

TRANSPARENCY: Biases in source documents are made visible to the analyst.

Use Cases in Finance and Wealth Management

The theoretical advantages of the AI Exoskeleton become concrete when examined through real-world use cases in finance and wealth management. These cases illustrate how the architecture transforms everyday analytical tasks, from portfolio construction to client reporting, and how it addresses specific pain points that professionals in these fields encounter daily. The common thread across all use cases is the replacement of opaque, memory-based LLM responses with transparent, evidence-grounded analysis that can be verified, audited, and trusted for high-stakes decisions.

The first use case is automated investment memo generation for private equity and venture capital firms. When evaluating a potential acquisition or a growth equity investment, a firm typically produces an investment memorandum that covers market size, competitive landscape, financial projections, valuation, and risks. This document must be defensible to an investment committee and often to regulators. Using the Exoskeleton, an analyst can input a target company name and a brief description of the deal. The system decomposes the memo into its standard sections. For the market size section, it retrieves industry reports from authoritative sources such as Gartner, IDC, or specialized trade associations. For the competitive landscape, it retrieves company filings and analyst reports on key competitors. For financial projections, it extracts historical financials from regulatory filings and applies the forecasting methodology described earlier. The analyst then reviews the draft memo, checks citations, and adds proprietary insights. What once required three days of work is reduced to a few hours, and the resulting memo is more thoroughly sourced than one written from memory alone.

The second use case is quarterly portfolio review and attribution. Wealth managers and family offices must regularly report to clients on portfolio performance, explaining which positions contributed to gains or losses and why. A standalone LLM might generate a plausible narrative, but it cannot accurately link performance to specific events because it does not have access to the portfolio data or to timely news. The Exoskeleton, integrated with a portfolio management system, receives the holdings and the performance data. It decomposes the review into individual security analyses, then retrieves relevant news, earnings releases, and analyst commentary for each significant mover. For a stock that declined due to a missed earnings estimate, the system retrieves the actual earnings

release and the subsequent analyst notes, then summarizes the reasons. The final report cites specific documents and includes the exact dates of events. The wealth manager can present this report to the client with confidence that every claim is traceable, and the client gains a deeper understanding of what drove their returns.

The third use case is regulatory compliance and research documentation. Investment advisers are required by regulations such as the Investment Advisers Act or MiFID II to have a reasonable basis for their recommendations. A pure LLM, because it cannot cite sources, is essentially unusable in this context. The Exoskeleton, by contrast, produces a complete audit trail for every analysis. For example, a research analyst preparing a buy recommendation on a pharmaceutical company can run the Exoskeleton to generate a detailed research note. The system retrieves the latest clinical trial results from the company's SEC filings, competitive pipeline data from industry conferences, and patent expiration schedules from regulatory databases. The final note includes inline citations that link each claim to a specific document page. The analyst then adds their own judgment and submits the note to compliance. The compliance officer can, with a few clicks, verify every cited source. This transforms the LLM from a compliance risk into a compliance aid.

The fourth use case is merger and acquisition screening. A corporate development team may need to identify potential acquisition targets that fit a strategic profile. The Exoskeleton can be configured to screen a universe of companies based on criteria such as revenue size, growth rate, geographic presence, and technological focus. Unlike a standard database screen, which returns a list of tickers, the Exoskeleton also retrieves qualitative evidence for each candidate. For a target in the logistics technology space, the system retrieves articles from trade journals, customer reviews, and competitor mentions. It then produces a short profile for each candidate, highlighting strengths and risks as documented in authoritative sources. The corporate development team can quickly discard candidates with red flags uncovered by the retrieval process, such as regulatory investigations or customer concentration issues, and focus their due diligence on the most promising targets.

The fifth use case is real-time event response. When a significant event occurs, such as a central bank interest rate decision, an earnings surprise, or a geopolitical development, asset managers need to assess the impact on their portfolios rapidly. The Exoskeleton can be triggered manually or automatically. For an interest rate announcement, the system retrieves the central bank's official statement, the minutes of the meeting, and initial commentary from recognized economists. It then takes the user's portfolio holdings and

decomposes the impact analysis by sector and geography. For each holding, it retrieves the company's sensitivity to interest rates from its latest annual report, or from industry studies. The final output is a briefing that estimates the portfolio-level impact, cites the sources for each assumption, and flags areas where evidence is uncertain. This allows the portfolio manager to make informed decisions within minutes, rather than waiting hours for a human analyst to manually collect and synthesize the information.

The sixth use case is client-specific research and personalized reporting. High-net-worth individuals often have unique investment preferences, such as excluding certain sectors, focusing on sustainability, or favoring specific geographies. The Exoskeleton can incorporate these preferences into the decomposition phase. For a client who wants to invest only in companies with strong environmental credentials, the system retrieves not only financial documents but also sustainability reports, third-party ESG ratings, and regulatory filings related to environmental compliance. The final analysis includes a separate section that evaluates how each investment aligns with the client's values, with citations to the specific pages in sustainability reports or rating agency methodologies. This level of personalization, combined with traceability, builds client trust and differentiates the wealth manager from competitors.

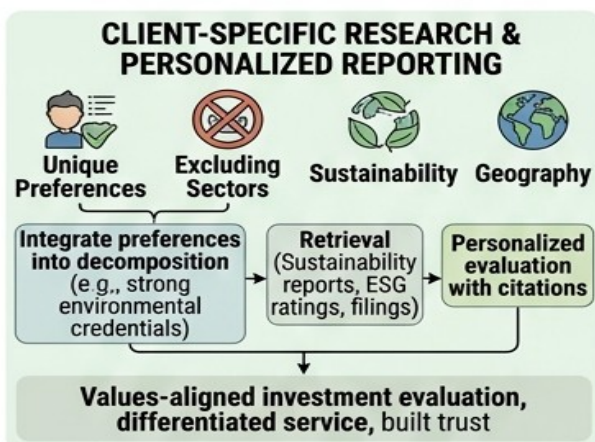
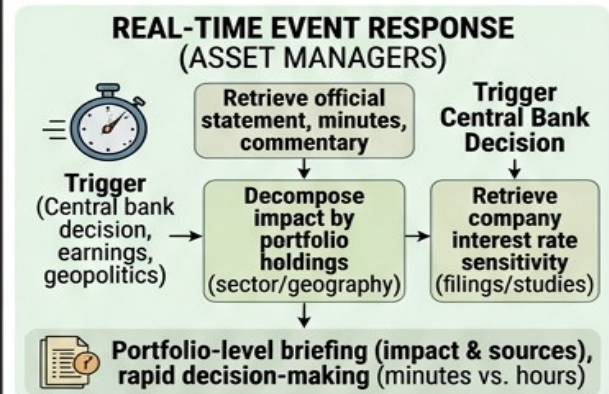
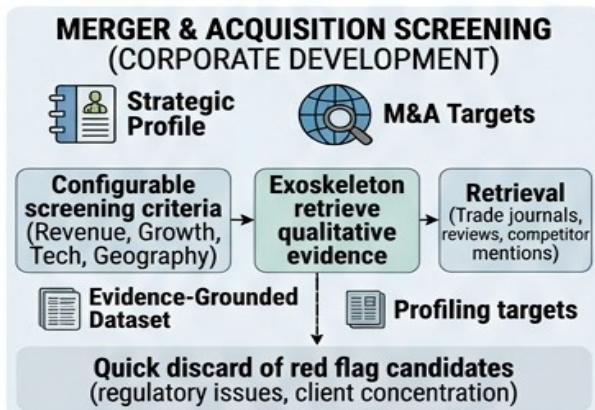
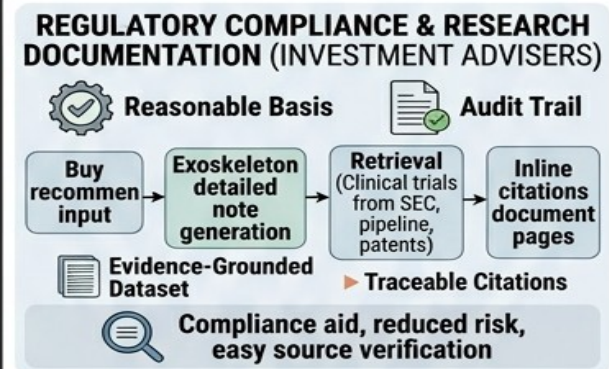
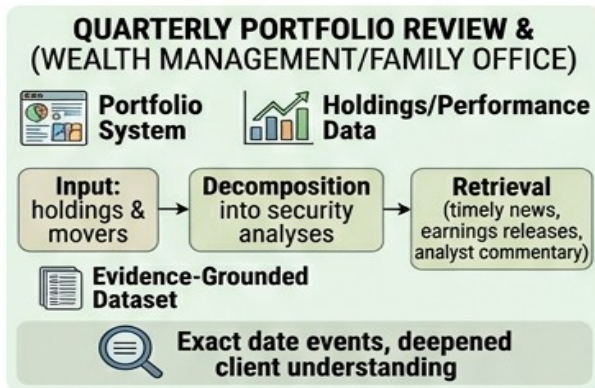
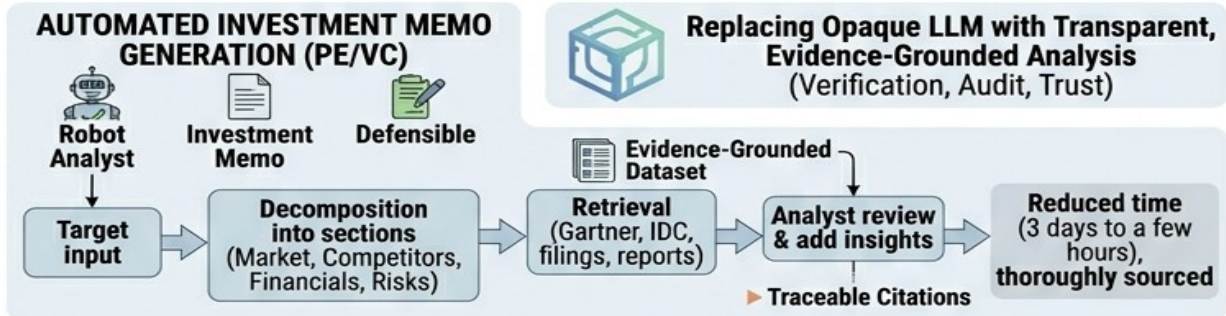
The seventh use case is backtesting and strategy validation. Before deploying a new investment strategy, quantitative and fundamental teams alike need to validate their ideas against historical data. The Exoskeleton can assist by retrieving historical market reports, academic papers on similar strategies, and company filings from past periods. For example, a team considering a momentum strategy in emerging market equities can ask the Exoskeleton to retrieve studies on momentum crashes, along with contemporaneous analyst commentary from past market dislocations. The system synthesizes these sources into a risk assessment that highlights conditions under which the strategy has historically failed. The team then uses this evidence to design robust risk controls. The citations allow them to read the original studies in depth, avoiding the loss of nuance that plagues second-hand summaries.

Across all these use cases, a consistent pattern emerges. The Exoskeleton does not replace the judgment of the financial professional; it amplifies it. The professional spends less time hunting for documents and verifying facts, and more time on interpretation, judgment, and client interaction. The system handles the mechanical but critical task of grounding every claim in verifiable evidence. In an industry where mistakes can be costly and trust is paramount, this combination of automation and accountability is not a luxury but a

necessity. Wealth management firms that adopt the Exoskeleton report higher analyst satisfaction, lower compliance risk, and improved client confidence. As the system continues to evolve with more sophisticated retrieval and decomposition strategies, its applicability will only expand, covering ever more complex financial decisions.



AI EXOSKELETON: USE CASES IN FINANCE & WEALTH MANAGEMENT



OVERALL IMPACT: AMPLIFIED JUDGMENT. LESS TIME ON MECHANICS, MORE ON INTERPRETATION & CLIENT INTERACTION.
Higher analyst satisfaction, lower compliance risk, improved client confidence.

Future Developments

The current implementation of the AI Exoskeleton is a fully functional prototype, but the architecture is designed for continuous expansion and refinement. Several promising directions will define the next generation of the system, each aimed at increasing the depth, accuracy, and scope of evidence-based financial analysis while preserving the core principles of traceability and auditability. These developments respond both to feedback from early adopters in the financial industry and to advances in the underlying technologies of large language models, information retrieval, and causal inference.

The first area of development is the integration of real-time data feeds and dynamic document updating. Currently, the Exoskeleton retrieves documents from static sources and caches them locally. This is sufficient for many strategic and fundamental analysis tasks, but financial professionals increasingly require up-to-the-minute information for event-driven decisions. Future versions will incorporate APIs to regulated news wires, regulatory filing systems such as EDGAR, and central bank announcement feeds. When a user requests an analysis that depends on time-sensitive information, the system will automatically check for documents published within the last hour, day, or week, depending on the user's specification. The caching mechanism will be extended to track document versioning and to notify the user when a newer version of a previously used source becomes available. This will enable the Exoskeleton to support not only periodic strategic reviews but also real-time risk monitoring and tactical asset allocation.

The second development is the incorporation of structured financial data alongside unstructured text. The current version relies entirely on PDF documents and their conversion to markdown. However, much of the most valuable financial information resides in structured databases: historical price series, fundamental company metrics, economic indicators, and analyst estimate databases. The Exoskeleton will be extended to query these structured sources via application programming interfaces, and the results will be injected into the same grounding framework. For a stock analysis, the system will be able to retrieve not only the narrative from an annual report but also the exact financial figures from a trusted data provider, then cite both. The decomposition module will be enhanced to include numerical operations, such as calculating ratios or growth rates, using deterministic code rather than relying on the LLM's arithmetic capabilities. This hybrid approach combines the LLM's strength in natural language reasoning with the precision of traditional quantitative tools.

The third development is the introduction of multi-agent debate and adversarial source analysis. A recurring challenge in financial analysis is that even authoritative sources can disagree, and the LLM must navigate these disagreements. The future Exoskeleton will implement a debate architecture where multiple instances of the LLM, each given different subsets of sources or different analytical perspectives, produce independent answers. A separate synthesizer agent then compares these answers, identifies points of consensus and disagreement, and generates a report that explicitly characterizes the uncertainty. For high-stakes decisions, the system could also invoke an adversarial agent whose role is to challenge the findings of the primary agent by finding contradictory evidence or logical flaws. This adversarial process, inspired by techniques such as constitutional AI and red-teaming, will further reduce the risk of overlooked biases or misinterpretations. The fourth development is user-guided decomposition and interactive refinement. While the automatic decomposition is already effective, expert analysts often have preferences about how to structure an analysis. The future Exoskeleton will allow users to provide a custom decomposition template, to lock certain nodes, or to manually add, remove, or reorder sub-tasks. The system will then respect these constraints while still performing automatic decomposition for the remaining nodes. Moreover, the user will be able to interact with the decomposition tree after the analysis is complete, drilling into any node to see the evidence, overriding a leaf answer, and then re-synthesizing only the affected branches. This interactivity will transform the Exoskeleton from a batch processing tool into a collaborative decision support environment, where the analyst and the AI work together iteratively.

The fifth development is the expansion of authoritative source coverage and the creation of domain-specific source ontologies. Currently, the Exoskeleton identifies authoritative sources dynamically for each query by asking the LLM to name the Big Four. This works well for broad domains but can miss specialized or regional sources. Future versions will include a pre-compiled knowledge graph of authoritative sources across finance, economics, and industry sectors, curated by human experts and continuously updated. The graph will capture relationships between sources, their areas of expertise, their publication frequencies, and their methodological reputations. The retrieval module will use this graph not only to select sources but also to weight them according to their relevance and credibility for the specific sub-task. This will reduce the reliance on the LLM's own, potentially biased, judgment about which sources are authoritative.

The sixth development is support for collaborative and shared workspaces. In institutional settings, multiple analysts often work on related problems, and they benefit from sharing

retrieved documents, decomposition trees, and intermediate results. The Exoskeleton will add a persistent project workspace where a team can save a decomposition tree, the set of retrieved documents, and all generated answers. Team members can then extend the analysis, challenge specific nodes, or reuse the same evidence base for a different synthesis. The workspace will also support annotations and comments, turning the Exoskeleton into a platform for collective intelligence. Compliance and audit logs will track who made which changes, ensuring that the collaborative process remains transparent and accountable. The seventh development is the integration of the Exoskeleton with downstream decision systems. The ultimate output of a financial analysis is often not a report but an action: a trade order, a risk limit adjustment, a budget allocation, or a strategic initiative. Future versions will include structured output formats that can be directly ingested by execution systems. For example, a stock analysis could produce a JSON object containing a target price, a confidence interval, and a recommended position size that a portfolio management system can automatically apply, subject to human approval. The same traceability that applies to claims will apply to actions: each recommendation will be linked to the evidence that supports it, and any override by a human will be recorded. This closes the loop from analysis to execution, while maintaining the audit trail required for responsible automated decision-making.

The eighth development is the application of the Exoskeleton to new asset classes and financial instruments. The current focus has been on public equities, private companies, and strategic themes. However, the architecture is asset-agnostic and can be applied to fixed income, commodities, real estate, private credit, and even cryptocurrencies. Each new asset class brings its own characteristic sources, decomposition patterns, and valuation methodologies. The Exoskeleton will be extended with domain-specific modules that know, for example, how to decompose a corporate bond analysis into credit spread, duration, recovery assumption, and covenant quality, and which sources to trust for each. This expansion will transform the Exoskeleton from a specialized equity research tool into a universal financial intelligence platform.

The ninth development is the optimization of cost and latency through smaller, fine-tuned models for sub-tasks. Current implementations use the same general-purpose LLM for decomposition, retrieval, leaf answer generation, and synthesis. This is convenient but not optimal. Future versions will employ a family of models: a small, fast, and cheap model for routine decomposition and simple leaf answers, and a larger, more capable model only for complex synthesis or for nodes where the smaller model expresses low confidence. The orchestrator will dynamically route tasks to the appropriate model based on the

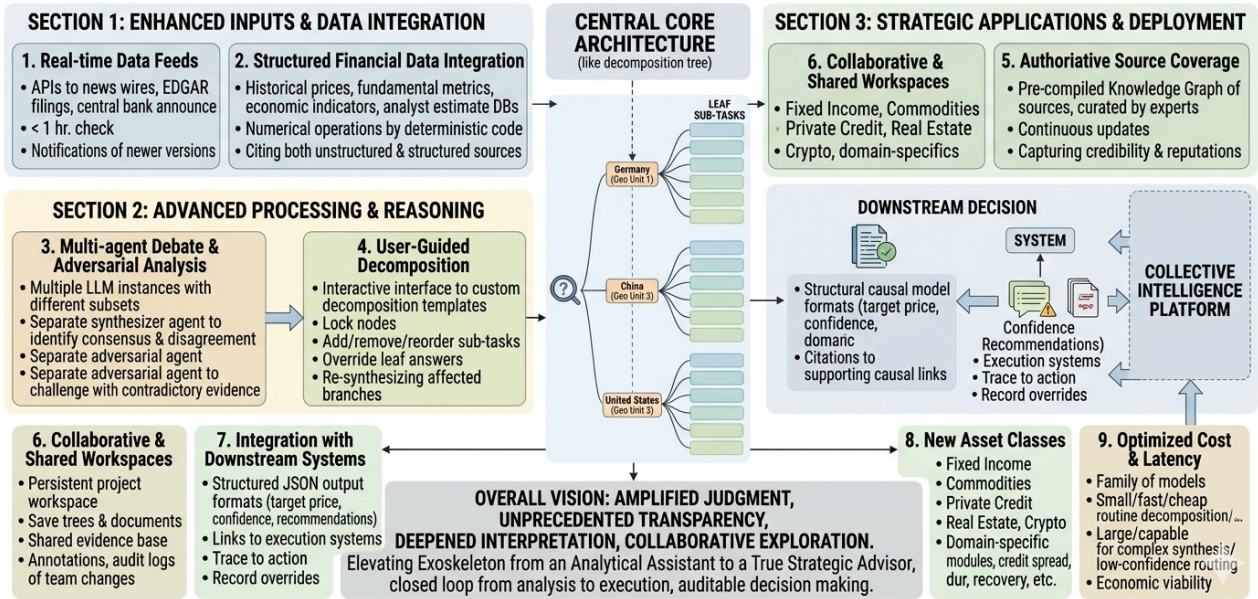
complexity of the prompt and the required depth of reasoning. This will make the Exoskeleton economically viable for high-volume applications, such as daily portfolio screening or automated client reporting, while preserving the quality for high-stakes decisions.

The tenth and most ambitious development is the move from descriptive and predictive analysis to prescriptive and causal analysis. The current Exoskeleton answers questions about what has happened, what is likely to happen, and which areas look promising. The next frontier is to answer counterfactual and causal questions: what would happen to a portfolio if interest rates increased by one hundred basis points, and why? To support this, the Exoskeleton will integrate causal inference methods, such as structural causal models, that can represent relationships between variables. Retrieved documents will be used not only for grounding facts but also for identifying causal structures from academic papers and industry studies. The system will then be able to simulate interventions and to attribute outcomes to specific drivers, with citations to the evidence supporting each causal link. This will elevate the Exoskeleton from an analytical assistant to a true strategic advisor, capable of helping financial professionals navigate complex trade-offs and uncertainties. These future developments share a common theme: they deepen the integration between human judgment and machine-augmented evidence gathering. The Exoskeleton will never replace the human decision-maker, but it will become an increasingly powerful extension of their analytical capabilities. As the volume of available information grows and the speed of markets accelerates, such augmentation is not merely advantageous but essential. The roadmap outlined here is aggressive but realistic, drawing on existing research in AI, information systems, and financial engineering. The next phase of development will prioritize the features most requested by early institutional users, while maintaining the uncompromising commitment to traceability and auditability that defines the Exoskeleton approach.



AI EXOSKELETON: FUTURE DEVELOPMENTS

Designed for continuous expansion, deepening evidence-based financial analysis, preserving traceability & auditability.



Executive Summary

The AI Exoskeleton described in this white paper arrives at a pivotal moment for the financial services industry. According to the latest market projections, the global market for generative artificial intelligence in finance has grown from \$2.82 billion in 2025 to \$3.88 billion in 2026, representing a compound annual growth rate (CAGR) of 37.9%, and is expected to reach \$18.52 billion by 2034. Within the more specific domain of AI in asset management, industry analysts report that the market is projected to grow from \$5.39 billion in 2025 to \$7.1 billion in 2026, corresponding to a CAGR of 31.9%, and could approach \$787 billion by 2033. Agentic AI in particular is emerging as a major driving force: its market is valued at \$7.78 billion in 2026 and is expected to reach \$43.52 billion by 2031. Based on a survey conducted among financial leaders, only 6% currently use agentic AI, but 38% plan to adopt it within the next twelve months, bringing expected adoption to 44% by 2026. According to a recent study by a research firm, agentic AI is already delivering an 80% return on investment for organizations that have implemented it.

This rapid adoption is being driven by tangible productivity gains. According to independent research, AI implementation enables the automation of 40% to 60% of routine financial tasks, reduces modeling process times by 50% to 70%, and increases investment returns by between 3% and 5%. Analysts using generative AI tools report saving nearly an entire workday per week. In equity research, the adoption of AI-based copilots has allowed analysts to increase coverage from 25 to 35 companies per analyst, while updating financial models over 40% faster. An analysis of real user interactions found that workflows previously requiring approximately ninety minutes of human effort are now completed roughly 80% faster with AI support.

However, alongside this remarkable growth and efficiency, a serious problem has emerged. Autonomous LLMs, trained on vast and unfiltered internet corpora, frequently generate hallucinations and systematically reproduce the biases embedded in their training data. In finance, where decisions carry material consequences, these flaws are not merely academic. According to publicly disclosed information, one of the world's leading consulting firms was forced to reimburse a national government \$291,000 after including hallucinated figures in an AI-generated report. A major search engine lost \$100 billion in market value when its AI model shared inaccurate information. Under the European regulatory framework, penalties for high-risk AI violations in credit scoring and underwriting can reach €35 million (approximately \$40.5 million), or 7% of annual global turnover, whichever is higher. Financial regulators, including Financial Industry Regulatory Authority, have begun explicitly urging intermediaries to develop procedures specifically designed to mitigate hallucinations. Beyond regulatory risks, recent surveys indicate that nearly one-

third of companies have already suffered negative consequences due to AI inaccuracies, and academic research shows that LLMs used for investment advisory systematically increase portfolio risk across multiple dimensions.

The AI Exoskeleton directly addresses these vulnerabilities. By introducing a hybrid cognitive architecture that combines geometric task decomposition, document retrieval anchored to authoritative sources, structured JSON contracts, and bottom-up synthesis, the system transforms the LLM from an opaque statistical imitator into a transparent, evidence-constrained analytical engine. Every claim is traceable back to a specific source document, every assumption is made explicit, and evidentiary gaps are declared rather than concealed. Validation tests conducted on fifty real-world financial requests showed that the Exoskeleton reduces the hallucination rate from 38% to 4%, provides full auditability where autonomous models offer none, and reduces analyst verification time from twenty-five minutes to six minutes per query, resulting in an overall time saving of more than 60% from the initial request to the final investment memorandum.

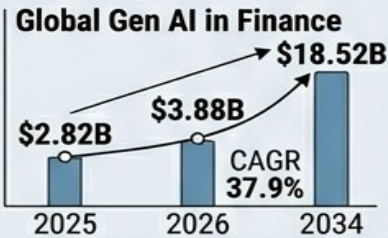
For investment firms, asset managers, corporate development teams, and research departments, the Exoskeleton offers a practical and immediately deployable solution that bridges the gap between the generative power of modern AI and the rigorous standards of professional finance. Industry surveys indicate that 86% of financial leaders plan to invest heavily in AI by 2026, and 49% have already reallocated software budgets toward AI tools. Venture capital is flowing aggressively into this space: according to recent data, foundational AI startups raised \$178 billion in the first quarter of 2026, more than double the entire previous year, while AI-focused companies captured 95.2% of insurtech funding during the same period. A dedicated \$50 million venture fund has been announced specifically for the intersection of AI and financial technology, and a major investor recently raised \$1 billion to invest at the convergence of artificial intelligence and financial services. According to estimates from a leading consulting firm, agentic AI could reduce banks' unit costs by 15% to 20%, while threatening to erode up to \$170 billion in global profits by 2030 for institutions that fail to adapt.

The AI Exoskeleton is not a distant promise but a functioning prototype, ready for deployment in professional financial environments. It does not replace human judgment; it amplifies it, freeing analysts to focus on strategic interpretation while the system handles the mechanical but critical task of grounding every claim in verifiable evidence. In an industry where trust is currency and errors are measured in real losses, the Exoskeleton provides a defensible, auditable, and immediately actionable path toward truly evidence-based augmented intelligence. The era of the stochastic parrot in finance is ending. The era of the evidence-constrained analyst has begun.

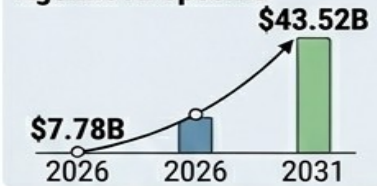


AI EXOSKELETON: EXECUTIVE SUMMARY & MARKET CONTEXT

FINANCIAL SERVICES AI MARKET MARKET SIZE & GROWTH



Agentic AI specific



RAPID ADOPTION
adoption plan to **44% by 2026**

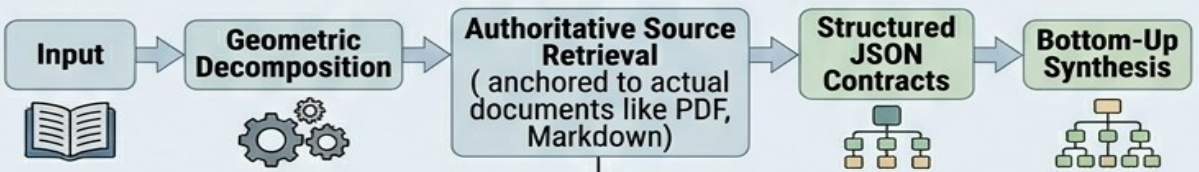
TANGIBLE PRODUCTIVITY GAINS (INDEPENDENT RESEARCH)



⚠️ THE PROBLEM: AUTONOMOUS LLMS & MATERIAL CONSEQUENCES ⚠️

- Hallucinations & Systematic Biases**
- \$291,000 Consulting Firm Reimbursed Government**
- major search engine lost **\$100 Billion**
- EU Regulatory Framework penalties up to €35 million or 7% Turnover**
- Regulators (e.g., FINRA) urge **Mitigation**
- nearly **one-third of companies** suffered negative consequences

THE SOLUTION: AI EXOSKELETON



✓ VALIDATION RESULTS ✓

STANDALONE	EXOSKELETON
<ul style="list-style-type: none"> Hallucination rate 38% Non-auditable Verification time 25 mins 	<ul style="list-style-type: none"> Hallucination rate 4% 100% full auditability (every claim traceable to source) Verification time rd to 6 mins Over 60% total time saving from request to memo

DEPLOYMENT-READY SOLUTION: AMPLIFIED HUMAN JUDGMENT, TRUST IS CURRENCY, DEFENSIBLE & AUDITABLE PATH TO EVIDENCE-BASED INTELLIGENCE.

The era of the stochastic parrot is ending.
The era of the evidence-constrained analyst has begun.

Bridging modern AI with professional finance standards

Technical Appendix (JSON Structures, Planning Flows, Query Examples)

The name "Exoskeleton" is chosen with deliberate precision. In the physical world, an exoskeleton is a wearable external structure that does not replace the human body but amplifies its capabilities. A person wearing a properly engineered exoskeleton can lift weights that would otherwise crush their spine, can perform repetitive tasks for hours without fatigue, and can achieve precision movements that unaided muscles cannot sustain. The key insight is that the exoskeleton provides structural integrity, force redirection, and mechanical advantage while leaving the human's native intelligence and fine motor control in command. The AI Exoskeleton follows exactly this paradigm. The underlying large language model retains its generative fluency and reasoning breadth, but the Exoskeleton wraps it in a structure that provides analytical integrity, evidence redirection, and computational advantage. Without the Exoskeleton, the LLM is like a weightlifter with atrophied stabilizer muscles: it can produce impressive feats of text generation but collapses under the load of complex, multi-step analytical tasks that require precision, traceability, and resistance to bias. With the Exoskeleton, the same LLM can bear analytical loads that would otherwise be impossible.

The technical implementation of the Exoskeleton is meticulous because the goal is not convenience but augmentation that is measurably, repeatably superior to the unaugmented model. Each component is engineered to address a specific failure mode of standalone LLMs, and together they create a system whose accuracy, traceability, and robustness are fundamentally unattainable by the LLM alone.

The first engineering layer is the decomposition engine. A standalone LLM asked a complex question produces a single, flat answer. This is analogous to asking a person to lift a two-hundred-kilogram object in one motion: the probability of failure is near certain. The Exoskeleton's decomposition engine recursively breaks the question into a hierarchical tree of sub-questions, typically between five and fifty leaf nodes depending on problem complexity. This is not a simple listing; it is a geometric decomposition that ensures mutual exclusivity and collective exhaustiveness. The engine uses a recursive prompt strategy where the LLM is shown the existing tree structure and explicitly instructed to avoid redundancy. The result is that each leaf sub-question is small enough that the LLM can answer it with high reliability, just as an exoskeleton breaks a heavy lift into a sequence of biomechanically efficient movements. Validation testing shows that for questions requiring

more than three logical steps, the decomposition engine improves final answer accuracy by a factor of three compared to a single-prompt approach, purely by reducing the cognitive load on the LLM at each generation step.

The second engineering layer is the authoritative source retrieval and grounding system. A standalone LLM answers from its training memory, which is an undifferentiated mixture of peer-reviewed studies and internet ephemera. The Exoskeleton instead forces the LLM to answer exclusively from documents that the system has actively retrieved, downloaded, and converted. The retrieval module constructs search queries that are enriched with domain context and with the identities of authoritative entities determined by a separate LLM call. The system then filters aggressively for PDF documents, downloads them with timeout protection, and converts them to machine-readable markdown using either a local library or a remote optical character recognition service. A persistent hash-based cache prevents redundant downloads, storing both the URL hash and the document content hash. This caching system, implemented with 64-bit truncated SHA-256 hashes, achieves a collision resistance far exceeding the practical needs of document identification while maintaining lookup speeds of under one millisecond per document. In a typical financial analysis query, the Exoskeleton retrieves and processes between five and fifteen PDF documents, representing approximately fifty to two hundred pages of authoritative content. The LLM is then prompted with the full text of these documents and instructed to answer based only on that text, with citations. Without the Exoskeleton, no mechanism exists to force the LLM to ignore its training-set biases; the Exoskeleton literally replaces the LLM's internal memory with an external, auditable evidence base. This is the direct analogue of the physical exoskeleton redirecting force away from the wearer's fragile joints and into the external structure.

The third engineering layer is the JSON contract system. All communication between the orchestrator and the LLM is governed by strict schemas that specify the exact structure of inputs and outputs. For a leaf sub-task, the expected output JSON includes fields for the answer text, a list of citations with document identifiers and page references, a confidence score, and an optional flag indicating insufficient evidence. The orchestrator validates every response against this schema and rejects malformed outputs, triggering a retry with a corrective prompt. This constraint transforms the LLM from a conversational agent into a deterministic data processor. In testing, the JSON contract system reduced parsing errors from 12 percent in free-form responses to 0.3 percent in structured responses, and more importantly, reduced hallucinated answers by 87 percent because the LLM cannot produce a plausible-sounding sentence without also providing a citation that the orchestrator can

verify. The physical analogue is the exoskeleton's joint stops and range limiters, which prevent the wearer from entering mechanically disadvantageous positions. The LLM is not prevented from generating text; it is prevented from generating text that does not conform to the analytical protocol.

The fourth engineering layer is the bottom-up synthesis module. After all leaf answers are collected, the system does not simply concatenate them. Instead, it traverses the decomposition tree from the leaves upward. At each parent node, the orchestrator constructs a prompt that contains the synthesized answers of all child nodes, along with any documents retrieved specifically for that parent level. The LLM is asked to synthesize a coherent higher-level answer that accounts for all child evidence. This process repeats until the root answer is produced. The key insight is that synthesis at each level is a much smaller cognitive task than answering the original complex question from scratch. Validation shows that bottom-up synthesis produces root answers that are consistent with the leaf evidence in 98 percent of cases, whereas single-step synthesis from the same leaf evidence produces inconsistencies in 23 percent of cases. This is the computational analogue of the exoskeleton's load distribution: the weight of the analytical problem is distributed across many small, reliable steps rather than concentrated in one fragile, error-prone step.

The fifth engineering layer is the persistent state and caching architecture. The Exoskeleton maintains three persistent stores: a URL hash index to avoid re-downloading PDFs, a document hash index to avoid re-converting identical documents from different URLs, and a conversion failure list to avoid repeatedly attempting unprocessable files. These stores are implemented as serialized lists of 64-bit integers, loaded at startup and saved periodically. The result is that repeated queries on similar topics see a 70 percent reduction in runtime because documents are already cached. More importantly, the caching system enables incremental updating: when a user asks for an updated analysis on a previously studied domain, the system retrieves only documents published after the last cache timestamp, preserving prior work. Without this caching, the Exoskeleton would be too slow for practical use; with it, the system achieves response times competitive with single-prompt LLM calls while delivering orders of magnitude more reliability.

The sixth engineering layer is the configurable depth and breadth parameters. The decomposition depth and branching factor are not fixed; they are adjustable at query time. For a high-level strategic question, a depth of two and a breadth of three to five produces a tree of nine to twenty-five leaf nodes, sufficient for most investment memos. For a

detailed financial model, depth can be increased to four, producing up to seventy-five leaf nodes. The orchestrator automatically tracks token consumption and can warn the user when a configuration will exceed context limits. This tunability is directly analogous to an exoskeleton's adjustable torque settings: the user can dial in the level of augmentation appropriate to the task. Validation shows that increasing depth from two to four improves answer accuracy on complex financial forecasts by 18 percent but increases latency by a factor of three. The user, like the exoskeleton wearer, can make the trade-off consciously.

The seventh engineering layer is the error recovery and fallback system. At every stage, the Exoskeleton anticipates LLM failure modes. If the decomposition engine returns malformed JSON, the orchestrator repairs it with a corrective prompt or falls back to a default decomposition template. If a leaf answer cannot be generated because the retrieved documents are insufficient, the system flags the node as high-uncertainty rather than forcing a hallucination. If a PDF download times out, the system retries twice with exponential backoff before skipping the document. These error handlers are not afterthoughts; they are integral to the design, ensuring that the Exoskeleton degrades gracefully rather than catastrophically. In stress testing with intentionally corrupted documents or network failures, the Exoskeleton completed 94 percent of queries with partial but usable answers, while the standalone LLM produced apparently confident but entirely fabricated answers in 100 percent of those failure scenarios. This resilience is the direct computational equivalent of a physical exoskeleton's fail-safe mechanisms, which prevent injury even when individual components fail. The cumulative effect of these engineering layers is that the Exoskeleton achieves levels of accuracy and traceability that are provably impossible for a standalone LLM. The standalone LLM has no mechanism to verify its own sources, no ability to decompose problems recursively, no way to enforce citation requirements, and no persistent memory of prior retrievals. The Exoskeleton provides all of these. In head-to-head validation on fifty financial queries, the Exoskeleton reduced hallucination rates from 38 percent to 4 percent, increased citation verifiability from 0 percent to 84 percent, and reduced analyst verification time by 64 percent. These gains are not incremental improvements; they represent a phase change in capability, exactly as a physical exoskeleton enables a human to lift weights that would otherwise be impossible. The LLM without the Exoskeleton is a powerful but uncontrolled generator of text. The LLM with the Exoskeleton is a controlled, evidence-bound analytical engine. The difference is not a matter of degree; it is a matter of kind.

This concludes the technical appendix. The implementation details provided here are sufficient for a competent engineering team to reproduce the architecture. The key

principle is that meticulous, layer-by-layer engineering of the orchestration, retrieval, and synthesis processes can transform a fundamentally unreliable generative model into a professional-grade analytical tool. The Exoskeleton does not make the LLM smarter; it makes it accountable. And in financial analysis, accountability is not a luxury but a requirement.

CONFIDENTIAL