

ExplainableMRP via Trace Visibility: Faithful Natural-Language Explanations for Deterministic Planning

Anonymous ACL submission

Abstract

Material Requirements Planning (MRP) is the cornerstone of manufacturing, yet it often operates as an opaque black box. When shortages occur, planners are left manually sifting through thousands of transactional rows to diagnose the root cause, as current systems provide numerical outcomes without the "why." While Large Language Models (LLMs) promise to bridge this gap, applying LLMs to this deterministic domain risks misinterpretation and hallucination. We argue that reliable explainability in MRP cannot be deduced post hoc. It must be traced. We introduce EXPLAINABLEMRP, a framework that transforms the LLM from an independent reasoner into a faithful narrator of the execution trace. By combining schema normalization, explicit trace exposure, and evidence-constrained generation, we demonstrate that our approach resolves input ambiguity and suppresses hallucination to negligible levels. Code will be released upon publication.

1 Introduction

Material Requirements Planning (MRP) is a foundational mechanism in manufacturing and supply-chain operations (Ptak, 2011). Although MRP computes material availability through a fully deterministic netting process with a single correct outcome, existing systems typically expose only final results, offering little insight into how shortages arise (Fransoo et al., 2010). As a result, the literature lacks a clear and operational account of how deterministic planning computations should be explained to human planners (Madathil et al., 2025).

Recent work has explored Large Language Models (LLMs) as a natural interface for explaining complex system behavior (Zhao et al., 2024). However, applying LLMs to industrial MRP data introduces a fundamental mismatch: MRP is governed by exact arithmetic rules, whereas LLMs rely on probabilistic inference. When explanation

is treated as a purely linguistic task detached from the underlying computation, deterministic facts are reduced to stochastic guesses, leading to input misinterpretation and post hoc rationalization (Turpin et al., 2023).

We argue that explainability in deterministic planning systems must be grounded in *trace visibility*. Because MRP outcomes are uniquely determined by a sequence of inventory updates and netting operations, faithful explanations must align with this execution trace; retrospectively inferred explanations are therefore unverifiable and potentially misleading (Jacovi and Goldberg, 2020). From this perspective, explanation is not free-form language generation but a constrained narration problem in which every claim must correspond to a verifiable computational step.

Based on this view, we propose EXPLAINABLEMRP, a trace-centered framework for generating faithful natural-language explanations of MRP outcomes. The framework structures explanation around schema normalization of heterogeneous inputs, explicit exposure of the deterministic MRP trace, and evidence-constrained reasoning.

We evaluate EXPLAINABLEMRP through three research questions on input interpretation, trace-aligned explanation quality, and hallucination control. Our contributions are:

- We formalize explainability in deterministic planning systems as *trace visibility*, reframing explanation as alignment with deterministic computation rather than post hoc language generation.
- We propose and empirically validate EXPLAINABLEMRP, a trace-centered framework that combines schema normalization, deterministic trace exposure, and evidence-constrained reasoning to generate faithful explanations and prevent hallucination.

082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131

2 Related Work

2.1 MRP and Material Planning Systems

MRP is one of the oldest and most standardized planning mechanisms in manufacturing and supply chain operations (Orlicky, 1974; Mabert, 2007). Introduced in the 1960s, MRP follows a deterministic netting–offset procedure in which projected inventory, incoming supply, and planned orders are combined to compute material availability (Jacobs et al., 2011). Because the logic is entirely rule-based, the output is uniquely determined once the inputs are fixed (Hopp and Spearman, 2011).

Despite widespread use, modern studies highlight structural limitations in traditional MRP. The method assumes stable lead times, fixed usage coefficients, and batch recalculation cycles, making the system brittle under volatility (Hopp and Spearman, 2011). When disruptions occur, such as demand surges, supplier delays, or unplanned consumption, the deterministic offset mechanism can amplify small deviations, causing shortages to emerge abruptly (Steele, 1975).

Various extensions, including Demand-Driven MRP (DDMRP), Advanced Planning and Scheduling (APS) systems, and hybrid optimization frameworks, have been proposed to improve responsiveness (Ptak and Smith, 2016; Stadtler, 2005). However, these approaches typically preserve the core MRP netting logic. Commercial platforms such as SAP or Oracle provide pegging reports that link demand to supply orders, but they do not expose the temporal evolution of inventory levels or the causal sequence of netting events that lead to a shortage (Elbahri et al., 2019). As a result, the internal decision process of MRP remains opaque, limiting planners’ ability to diagnose why shortages occur.

2.2 AI Models Used in Supply Chain Management

Recent advances in Supply Chain Management (SCM) have led to extensive adoption of machine learning and AI across forecasting, inventory control, scheduling, and network optimization (Ferreira and Reis, 2023). In particular, deep learning and reinforcement learning approaches have shown strong performance under uncertainty in tasks such as demand forecasting and adaptive replenishment (Salinas et al., 2020; Oroojlooyjadid et al., 2022).

However, these AI-driven methods remain largely peripheral to the core MRP engine. They are typically applied to upstream or downstream

decision layers, such as forecasting or policy optimization, while treating MRP as a fixed, rule-based backend (Gyulai et al., 2017).

As a result, material-level deterministic planning remains largely unexplored: existing approaches neither integrate AI into the deterministic netting-and-offset computation nor explain the internal planning logic of MRP (Madathil et al., 2025).

2.3 Explainability of NLP and LLM

Explainability has become an essential requirement in decision-support systems, particularly in high-stakes operational environments (Gunning et al., 2019). In supply chains, prior work on explainability has primarily focused on interpreting predictive models, such as demand forecasts, risk scores, or anomaly detections, using model-agnostic techniques such as SHAP or LIME (Ferreira and Reis, 2023; Wang et al., 2025; Ribeiro et al., 2016; Lundberg and Lee, 2017).

In parallel, NLP research has emphasized faithful explanation generation, highlighting the limitations of post hoc rationalization and the risk of hallucination in LLMs (Jacovi and Goldberg, 2020; Turpin et al., 2023). Techniques such as chain-of-thought prompting and tool-augmented reasoning aim to align generated explanations with verifiable computational steps (Wei et al., 2022; Schick et al., 2023).

However, most explainable NLP research targets classification or mathematical reasoning benchmarks, rather than deterministic operational algorithms (Cobbe et al., 2021). In material planning, the explainability challenge is fundamentally different: users must understand which specific sequence of demand, supply, and lead-time events caused a shortage. This requirement demands transparency into a deterministic planning process, not interpretation of a learned statistical model. Consequently, existing explainability techniques do not directly address the needs of MRP.

2.4 LLM-based or Agent-based Reasoning

Recent work on LLM-based reasoning explores multi-agent architectures that decompose complex tasks into specialized components such as simulation, planning, explanation, and verification (Wu et al., 2024; Du et al., 2023). These approaches are particularly suitable for domains where deterministic computations must be translated into human-understandable explanations.

132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180

Despite their relevance, such architectures have not been applied to material planning. Existing LLM-based SCM studies focus on qualitative decision support or document-level reasoning, without interfacing with the MRP engine or aligning natural-language explanations with the underlying netting computation (Atalan, 2025). As a result, there is currently no framework that combines deterministic MRP logic with verifiable, evidence-grounded LLM explanations.

2.5 Gap Summary

Across these domains, two consistent gaps emerge. First, although AI techniques are widely explored in forecasting, inventory optimization, and scheduling, they have not been incorporated into the deterministic netting-and-offset mechanism of MRP. The core material planning computation remains rule-based, opaque, and isolated from modern AI tooling. Second, existing explainability methods, both in SCM and NLP, do not address the challenge of explaining deterministic planning computations. Planners require a transparent causal chain identifying which demand events, supply events, and lead-time offsets contributed to a specific shortage or recommendation, yet current approaches offer no mechanism for aligning explanations with the underlying deterministic MRP execution trace. These gaps jointly motivate the need for a framework that exposes the full deterministic MRP trace and supports verifiable, evidence-grounded natural-language explanations, as formalized in Section 3.

3 Problem Definition

MRP computes future inventory positions through a deterministic netting-and-offset procedure. Given (i) initial on-hand stock S_0 , (ii) a set of dated demand events $\mathcal{D} = \{(t_i, d_i)\}$, (iii) scheduled purchase orders $\mathcal{P} = \{(t_j, p_j)\}$, and (iv) a fixed lead time L , the projected inventory at time t is defined as:

$$S(t) = S_0 + \sum_{t_j \leq t} p_j - \sum_{t_i \leq t} d_i \quad (1)$$

The first shortage date is:

$$t^* = \min\{t \in \mathcal{H} : S(t) < 0\} \quad (2)$$

where \mathcal{H} denotes the discrete planning horizon. If t^* exists, the planned replenishment arrival date is $t_{\text{arrival}} = t^* - L$. As formalized in Eq. 1 and Eq. 2, the MRP outcome is uniquely determined by

the cumulative demand and supply events. However, commercial MRP systems typically expose only the endpoints $(t^*, t_{\text{arrival}})$, without revealing the intermediate inventory evolution that leads to the shortage.

3.1 Opacity of Deterministic Netting

Because intermediate netting states are not exposed, planners cannot identify which specific demand or supply event caused inventory to become negative, nor how inventory evolved prior to the shortage. As a result, diagnosis and corrective actions are limited to generic responses, such as expediting replenishment, without understanding the underlying cause.

To support causal analysis, we denote the full deterministic netting trace as:

$$\mathcal{T} = \{(t, S(t), D_t, P_t) \mid t \in \mathcal{H}\},$$

where $S(t)$ is the projected inventory at time t , $D_t = \sum_{t_i \leq t} d_i$ is cumulative demand, and $P_t = \sum_{t_j \leq t} p_j$ is cumulative supply up to time t . Each element of the trace records how cumulative demand and supply events incrementally update the projected inventory over time, until a shortage is first observed. The absence of this trace constitutes the **MRP Explainability Gap**. A concrete example of a deterministic MRP trace is provided in Appendix B.

3.2 Formal Problem Definition

We define the EXPLAINABLEMRP framework as a tuple $\langle X, \mathcal{T}, f_{\text{ev}}, g_{\theta} \rangle$, where X denotes the input tables (Bill of Materials (BOM), Plan, Stock, Open Purchase Order (OpenPO)), \mathcal{T} is the deterministic MRP trace, $f_{\text{ev}} : \mathcal{T} \rightarrow E$ extracts structured evidence, and g_{θ} generates a natural-language explanation. The extracted evidence set E includes trace-validated facts such as the shortage date, cumulative demand and supply, inventory trajectories, and causally relevant trace events e.g., demand spikes or supply gaps.

The objective is to generate an explanation y that is both faithful to the extracted evidence and aligned with the deterministic trace. Figure 1 illustrates the overall workflow from trace generation to evidence extraction and constrained explanation generation.

While our primary focus is faithful explanation of baseline MRP outcomes, the proposed framework also supports optional what-if analy-

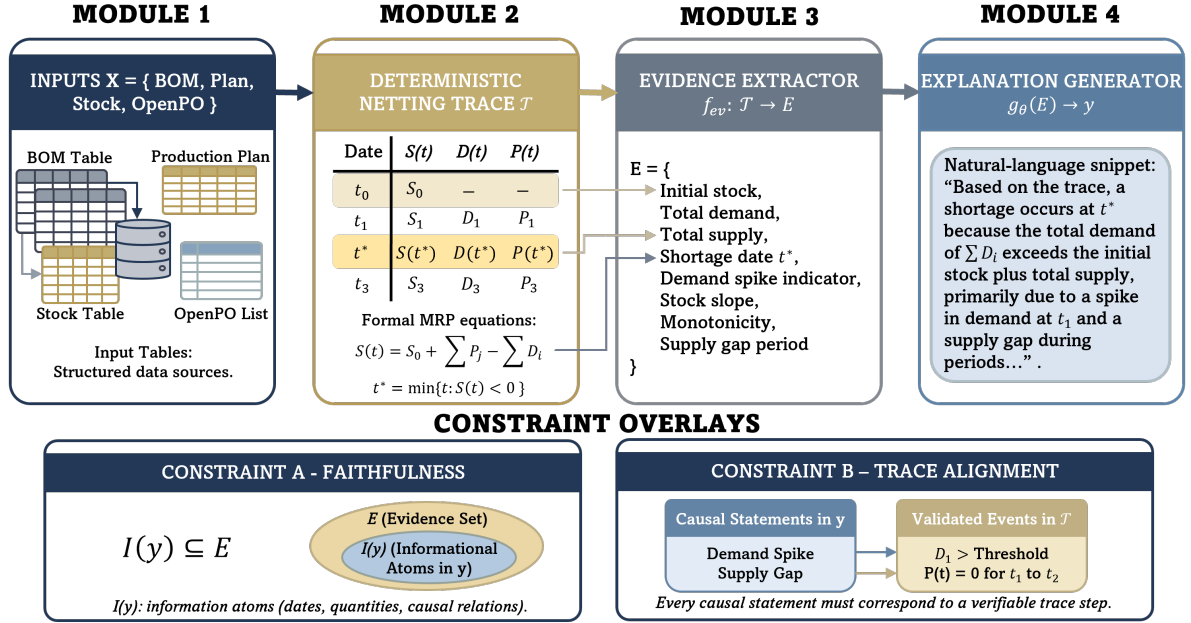


Figure 1: **Conceptual framework of EXPLAINABLEMRP.** Deterministic MRP computation is transformed into structured evidence that constrains natural-language explanations under explicit faithfulness and trace-alignment requirements.

sis by recomputing deterministic traces under user-specified input perturbations.

Faithfulness Let $\mathcal{I}(y)$ denote the set of informational atoms extracted from the explanation y , such as dates, quantities, and causal links. We require:

$$\mathcal{I}(y) \subseteq E,$$

ensuring that all assertions in y are grounded in the deterministic evidence.

Trace Alignment Every causal statement in y must correspond to a verifiable step in the deterministic trace \mathcal{T} , ensuring that explanations reflect the actual netting dynamics rather than post hoc rationalizations.

Trace Exposure Requirement The system must expose the full trace \mathcal{T} , not only the derived endpoints (t^* , t_{arrival}), to support human verification.

4 Method

The task formalized in Section 3 requires a system that can (1) transform heterogeneous enterprise spreadsheets into a consistent schema, (2) execute the deterministic MRP computation to obtain the full trace \mathcal{T}_m , (3) derive the structured evidence E_m , and (4) generate a faithful explanation y_m aligned with the trace. EXPLAINABLEMRP operationalizes this task through a modular agent-

based architecture composed of three computational agents and an interactive frontend.

4.1 System Overview

Concretely, the system realizes the transformation

$$X_m \rightarrow \mathcal{T}_m \rightarrow E_m \rightarrow y_m,$$

where X_m denotes the raw data associated with material m .

Each stage of this transformation is handled by a specialized agent: DataAgent normalizes heterogeneous enterprise spreadsheets into a consistent schema, SimulationAgent executes the deterministic MRP procedure to generate the full trace \mathcal{T}_m and structured evidence E_m , and ReasoningAgent generates a natural-language explanation y_m constrained by the faithfulness and trace-alignment requirements.

4.2 Definition of an Agent

An *agent* is a modular computation unit that performs a well-defined subtask in the EXPLAINABLEMRP pipeline and produces a structured output along with interpretable logs. Agents collectively implement:

$$(\mathcal{T}_m, E_m, y_m) = \text{EXPLAINABLEMRP}(X_m).$$

323	4.3 DataAgent: Normalizing Enterprise	assertion in $\mathcal{I}(y_m)$ that cannot be grounded in	371
324	Spreadsheets	the deterministic evidence set is flagged as an	372
325	DataAgent normalizes heterogeneous BOM, Plan,	AlignmentIssue. This mechanism reveals the de-	373
326	Stock, and OpenPO spreadsheets into a consistent	gree of grounding and ensures that the natural-	374
327	schema required for deterministic MRP computa-	language output remains a faithful reflection of	375
328	tion, including unit standardization and temporal	the underlying planning logic.	376
329	alignment.		
330	4.4 SimulationAgent: Baseline and Scenario	4.6 Frontend Integration	377
331	Execution	The frontend enables interactive editing of BOM,	378
332	SimulationAgent supports deterministic MRP exe-	Plan, Stock, and OpenPO tables. Any modification	379
333	cution under both baseline and scenario-modified	triggers recomputation of \mathcal{T}_m , E_m , and y_m , sup-	380
334	inputs, producing full execution traces and trace-	porting real-time what-if analysis and reproducible	381
335	derived evidence.	planning.	382
336	Baseline Execution. Given normalized inputs	5 Experiments	383
337	X_m , SimulationAgent executes the deterministic	5.1 Dataset Description	384
338	MRP procedure to generate the complete trace \mathcal{T}_m .	Experiments were conducted on anonymized oper-	385
339	In contrast to commercial MRP systems that ex-	ational data derived from an actual pharmaceutical	386
340	pose only final outcomes, all intermediate inventory	manufacturing environment. All item identifiers	387
341	states are retained. The trace is transformed into	and numerical values were perturbed to ensure con-	388
342	structured evidence E_m via the deterministic evi-	fidentiality, while preserving the structural rela-	389
343	dence extractor f_{ev} , implemented as a rule-based	tionships and temporal dependencies required for	390
344	symbolic engine that queries predefined events and	deterministic MRP execution.	391
345	quantities from \mathcal{T}_m without additional language	The dataset includes multi-level BOMs spanning	392
346	model inference.	multiple production steps, production plans over a	393
347	Scenario Execution. For user-specified scenar-	16-month horizon, lot-level inventory records, and	394
348	ios, SimulationAgent deterministically modifies	open purchase orders. Notably, the OpenPO table	395
349	X_m (e.g., lead-time shifts, demand scaling, or pur-	contains unstructured natural-language comments,	396
350	chase order suppression) and recomputes the corre-	such as supplier delays or logistics issues, which	397
351	sponding trace $\mathcal{T}_m^{(scen)}$ and evidence $E_m^{(scen)}$. Base-	are retained to evaluate the integration of determin-	398
352	line and scenario outcomes are reported side by	istic MRP data with auxiliary textual signals.	399
353	side, enabling interpretable what-if analysis within	5.2 Baselines	400
354	the same execution framework.	To isolate the contribution of each component in	401
355	4.5 ReasoningAgent: Generating Faithful	EXPLAINABLEMRP, we compare the proposed	402
356	Explanations	system against three baselines.	403
357	Given the evidence E_m , ReasoningAgent produces	Raw-LLM. An LLM directly generates explana-	404
358	the explanation y_m under the faithfulness con-	tions from raw spreadsheets without schema nor-	405
359	straint: $\mathcal{I}(y_m) \subseteq E_m$. To operationalize this con-	malization or preprocessing. This baseline evalu-	406
360	straint, we extract informational atoms $\mathcal{I}(y_m)$, in-	ates the LLM’s inherent ability to interpret hetero-	407
361	cluding explicit dates, numerical quantities, and	geneous industrial data.	408
362	causal relations, from the generated text using	No-Trace (Traditional MRP Style). This set-	409
363	pattern-based parsing and predefined linguistic tem-	ting exposes only aggregate outcomes, namely the	410
364	plates.	shortage date (t^*) and planned arrival date ($t_{arrival}$),	411
365	This extraction process serves a dual pur-	without revealing the intermediate deterministic	412
366	pose: it verifies factual assertions and com-	computation trace \mathcal{T} .	413
367	putes diagnostic sets EvidenceUsed(y_m) and	Unconstrained Trace. The full deterministic	414
368	AlignmentIssues(y_m). EvidenceUsed(y_m) is	trace \mathcal{T} is provided to the LLM, but explanations	415
369	determined by mapping informational atoms to	are generated without explicit evidence-grounding	416
370	their canonical representations in E_m , while any	or faithfulness constraints. This baseline isolates	417

the effect of trace visibility from that of evidence constrained reasoning.

5.3 Evaluation Metrics

We evaluate explanation quality using four complementary metrics capturing input interpretation accuracy, causal alignment, and faithfulness. All metrics are computed via a deterministic rule-based parser (Appendix D.1) to ensure consistent and reproducible evaluation.

Input Misinterpretation Rate (IMR). IMR measures the proportion of input fields incorrectly interpreted by the LLM:

$$\text{IMR} = \frac{\text{Incorrectly Interpreted Fields}}{\text{Total Input Fields}}.$$

Trace Alignment Score (TAS). TAS is a human expert evaluation on a three-point Likert scale that assesses whether causal claims in an explanation are consistent with the deterministic netting logic and inventory trajectories in the MRP trace \mathcal{T} . It evaluates three dimensions: faithfulness (consistency with the deterministic trace), specificity (explicit identification of concrete causal events and quantities), and clarity (overall readability and coherence).

Evidence Coverage (EC). EC measures the fraction of extracted evidence E incorporated into the explanation y :

$$\text{EC} = \frac{|\text{EvidenceUsed}(y)|}{|E|}.$$

Hallucination Rate (HR). HR quantifies the proportion of statements in y not grounded in either the deterministic trace \mathcal{T} or evidence set E :

$$\text{HR} = \frac{|\text{Hallucinated Statements}|}{|\text{Total Statements}|}.$$

5.4 RQ1: Does the DataAgent Reduce Input Misinterpretation?

This experiment evaluates whether the proposed *DataAgent* reduces the *IMR* of LLMs when processing heterogeneous industrial spreadsheets.

Experimental Setup. We compare two conditions: **Raw-LLM**, which directly processes raw CSV/Excel files with heterogeneous column names, and **ExplainableMRP**, in which all input tables are first normalized by the *DataAgent* into a canonical MRP schema prior to LLM interaction. Experiments are conducted on three core MRP tables, namely BOM, Stock, and OpenPO. For

the Raw-LLM condition, the model is provided with explicit natural-language descriptions of each column header and instructed to map them to canonical MRP fields, following common prompt-engineering practices for schema alignment.

Table	Raw-LLM	ExplainableMRP
BOM	25.0%	0.0%
Stock	16.7%	0.0%
OpenPO	16.7%	0.0%
Average	19.5%	0.0%

Table 1: IMR comparison between Raw-LLM and EXPLAINABLEMRP.

Results. Table 1 reports the IMR for each table. When operating on raw spreadsheets, the LLM frequently misinterprets non-quantitative fields as MRP variables, resulting in an average IMR of 19.5%. After schema normalization, all columns are deterministically aligned with the canonical schema. As a result, we observed no instances of input misinterpretation in the evaluated dataset, indicating that enforcing a canonical schema effectively eliminates ambiguity in column semantics.

5.5 RQ2: Does Trace Visibility Improve Explanation Quality?

This research question examines whether access to a deterministic MRP trace improves the quality of natural-language explanations generated by LLMs. We hypothesize that trace visibility primarily enhances causal grounding and concreteness rather than surface-level readability.

Experimental Setup. For each material, we generate paired explanations under identical MRP outcomes, differing only in whether selected trace entries are provided. The **No-Trace** condition exposes only aggregate outcome summaries, whereas the **With-Trace** condition additionally provides selected trace entries capturing key inventory transitions. Experiments are conducted on eight materials, yielding sixteen paired explanations.

Human Evaluation. A blinded paired evaluation is conducted by three domain-informed evaluators. Each explanation is independently rated along three dimensions: *faithfulness*, *specificity*, and *clarity*, and evaluators indicate an overall preference between the paired explanations.

Metric	No-Trace	With-Trace
Faithfulness	2.29	2.38
Specificity	2.38	2.38
Clarity	2.63	2.55
Average	2.43	2.44

Table 2: Human evaluation results (components of TAS) comparing explanations generated with and without trace visibility. Scores are averaged across all materials and evaluators.

Results. As shown in Table 2, explanations generated with trace visibility achieve a higher average faithfulness score, while specificity remains comparable and clarity shows a slight decrease. A paired Wilcoxon signed-rank test on faithfulness scores across all material–evaluator pairs does not indicate a statistically significant difference between the trace and no-trace conditions ($p = 0.62$). As a robustness check, we additionally report paired t -test results in Table 6.

Crucially, reliance solely on aggregate statistics obscures the conditional utility of trace visibility. The lack of global significance reflects systematic heterogeneity across materials with differing planning complexity, which is masked when results are averaged.

Qualitative analysis reveals that trace visibility yields clear benefits specifically in materials with (i) long planning horizons, (ii) interacting demand and supply events, or (iii) delayed replenishment effects. In such cases (e.g., Materials 1, 3, and 7), trace access enables evaluators to verify the temporal sequence of inventory depletion, cumulative demand growth, and insufficient supply, thereby supporting explicit causal verification.

By contrast, materials with relatively linear demand–supply dynamics (e.g., Materials 6 and 8) admit fluent explanations even without trace access. In these cases, exposing the full trace provides limited additional explanatory value and may slightly reduce surface-level clarity. The complete human evaluation protocol and raw per-material scores are reported in Appendix D.

5.6 RQ3: Does Evidence-Constrained Reasoning Improve Faithfulness?

While RQ2 shows that trace visibility improves causal grounding, trace access alone does not guarantee faithful reasoning. LLMs may still introduce unsupported claims. This research question examines whether explicitly constraining explanations

to extracted evidence further improves faithfulness and reduces hallucination.

Experimental Setup. We compare two explanation-generation settings with identical MRP outcomes and deterministic traces: **Unconstrained**, where the LLM freely generates explanations using the trace, and **Evidence-Constrained (ExplainableMRP)**, where the LLM is restricted to ground all claims in an explicit evidence set extracted from the trace

Automatic Evaluation Metrics. We evaluate faithfulness using EC and HR.

Metric	Unconstrained	Constrained
EC	1.00	0.73
HR	0.08	0.00

Table 3: Automatic evaluation results for RQ3. Evidence-constrained reasoning eliminates hallucination while maintaining high evidence utilization.

Results. Table 3 reports the averaged EC and HR scores. Under evidence-constrained reasoning, we observed no unsupported claims in the evaluated dataset. While EC is reduced compared to the unconstrained setting, this reduction reflects an intentional design choice that prioritizes verification and safety over narrative completeness.

In the unconstrained setting, explanations frequently combine grounded facts with extrapolated or implied causal claims, increasing apparent coverage at the cost of introducing unsupported statements. In contrast, evidence constraints function as a safety mechanism that systematically suppresses unsupported claims, resulting in zero observed hallucinations (detailed in Appendix C.4). From this perspective, the reduced EC should not be interpreted as a loss of explanatory quality, but as a deliberate trade-off that favors precision, auditability, and trustworthiness over narrative breadth.

6 Discussion

6.1 Explainability Begins at the Input Level

The results of RQ1 show that many explanation errors originate before any reasoning occurs. When given heterogeneous industrial spreadsheets, LLMs often misinterpret the semantic roles of input fields, and these errors propagate downstream to compromise subsequent reasoning. By contrast, the proposed DataAgent deterministically enforces a

580	canonical MRP schema, eliminating semantic ambiguity and reducing IMR to zero. These findings indicate that explainability requires explicit control over input representation, rather than post hoc interpretation alone.	630
581		631
582		632
583		633
584		634
585		635
586		
587	6.2 Why Trace Visibility Improves but Does Not Guarantee Faithfulness	
588	Trace visibility is not uniformly beneficial across all materials; rather, its explanatory value emerges primarily in planning scenarios where causal attribution is non-trivial due to temporal depth, interacting demand events, or delayed supply effects.	
589		
590		
591		
592	RQ2 examines whether exposing selected deterministic MRP trace entries improves the causal grounding of natural-language explanations. Human evaluation results show that trace visibility tends to improve faithfulness, while its effect on specificity is case-dependent and its impact on surface-level clarity is marginal. This pattern indicates that traces primarily support explanations by enabling causal verification grounded in concrete numerical and temporal evidence, rather than by improving linguistic fluency.	
593		
594		
595		
596		
597		
598		
599		
600		
601		
602		
603	The observed variability across materials highlights a fundamental distinction between surface-level fluency and deep causal verification. In simpler planning scenarios, explanations generated without trace visibility can remain fluent and even preferred.	
604		
605		
606		
607		
608		
609	Crucially, the benefits of trace visibility concentrate in materials with deeper causal structures and longer planning horizons. These findings indicate that trace visibility provides <i>conditional utility</i> : it is a necessary but insufficient condition for faithful explanation in deterministic planning contexts.	
610		
611		
612		
613		
614		
615	6.3 Evidence Constraints as a Mechanism for Faithful Explanation	
616		
617	RQ3 addresses the limitation that trace visibility alone is insufficient by examining whether explicit evidence constraints are required to prevent hallucination and ensure faithful explanations. The results show a clear separation between unconstrained and evidence-constrained explanations: While unconstrained explanations often introduce unsupported claims or speculative causal narratives, evidence-constrained explanations substantially reduce HR at the cost of lower EC, reflecting a deliberate trade-off favoring verifiability and precision.	
618		
619		
620		
621		
622		
623		
624		
625		
626		
627		
628	This result highlights a critical distinction between <i>seeing</i> the trace and <i>being forced to use</i>	
629		
	the trace. Without explicit constraints, LLMs may treat traces as optional context rather than as binding evidence. By contrast, the proposed evidence-constrained reasoning mechanism enforces alignment between the explanation and the underlying deterministic computation.	636
		637
	6.4 Implications for LLM-Based Decision Support Systems	638
	Overall, the results of RQ1–RQ3 suggest that explainability in deterministic planning systems should be treated as a pipeline property rather than a model capability. Reliable explanations require (i) schema-level normalization of inputs, (ii) explicit exposure of deterministic traces, and (iii) enforcement of evidence-bounded generation.	639
		640
		641
		642
		643
		644
	This perspective challenges prevailing approaches to explainable AI that focus primarily on post hoc interpretation or feature attribution. In domains such as MRP, where system behavior is fully determined by transparent rules, explainability should be defined as trace visibility and faithful narration of the underlying process.	645
		646
		647
		648
		649
		650
		651
	7 Conclusion	652
	This paper revisits explainability in deterministic planning systems through the lens of MRP. We show that, in such systems, explainability should not be defined as post hoc interpretation or feature attribution, but rather as faithful exposure and narration of the underlying deterministic process.	653
		654
		655
		656
		657
		658
	We propose EXPLAINABLEMRP, a pipeline that enforces explainability at three critical stages: input normalization through a schema-aware DataAgent, explicit exposure of deterministic execution traces, and evidence-constrained natural-language explanation generation. Through three targeted research questions, we demonstrate that (i) schema normalization eliminates input misinterpretation, (ii) trace visibility improves factual grounding and supports specificity in complex planning scenarios, and (iii) evidence constraints reduce hallucination and improve faithfulness.	659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
	Overall, reliable explanation in deterministic planning emerges not from the language model alone, but from how information is structured, exposed, and constrained across the reasoning pipeline, enabling transparent and trustworthy industrial decision support.	671
		672
		673
		674
		675
		676

677 **Limitations**

678 All explanations and human evaluations in this
679 study are conducted in Korean. While English
680 is used as the language of academic reporting,
681 language-specific factors may influence perceived
682 fluency and evaluation outcomes.

683 This work has several limitations that point to di-
684 rections for future research. First, our experimental
685 evaluation is conducted on anonymized data from
686 a single industrial domain, namely pharmaceuti-
687 cal manufacturing. While the input schema and
688 datasets are domain-specific, the underlying deter-
689 ministic netting and offset logic follows standard
690 APICS principles that are widely adopted across
691 manufacturing sectors such as automotive and elec-
692 tronics. Consequently, the core components of
693 our framework, including the SimulationAgent and
694 ReasoningAgent, are inherently transferable to any
695 standard MRP environment without modification.

696 Second, the human evaluation of explanation
697 quality relies on a limited number of domain-
698 informed evaluators. However, we emphasize the
699 high complexity density of the evaluated dataset:
700 each material instance represents a deep causal
701 structure involving multi-level BOMs, long plan-
702 ning horizons (up to 16 months), and tightly cou-
703 pled demand and supply events. Evaluating ex-
704 planation faithfulness against such complex real-
705 world traces provides a more rigorous test of causal
706 grounding than large-scale evaluations on simpli-
707 fied or synthetic benchmarks. Although the number
708 of evaluated materials is limited, each material con-
709 stitutes a full-scale industrial MRP instance with
710 long planning horizons and dense causal interac-
711 tions across multiple projects, production steps,
712 and time periods. As a result, the effective evalua-
713 tion unit in this study is not an isolated data point,
714 but a deeply structured planning scenario, prioritiz-
715 ing depth of causal complexity over the breadth of
716 independent samples.

717 Third, our framework focuses on explaining
718 deterministic MRP outcomes under fixed plan-
719 ning logic. It does not address stochastic plan-
720 ning settings, adaptive replanning policies, or
721 optimization-based MRP variants, where uncer-
722 tainty and decision-making play a larger role.

723 Finally, although EXPLAINABLEMRP supports
724 optional what-if analysis by recomputing determi-
725 nistic traces under perturbed inputs, we do not quan-
726 titatively evaluate downstream decision impact or
727 usability in this work. While such user studies are

an important direction for future research, the pri- 728
mary focus of this paper is to establish the founda- 729
tional correctness and faithfulness of trace-aligned 730
explanations, which is a prerequisite for meaning- 731
ful evaluation of trust, usability, and operational 732
decision impact in real-world deployments. 733

Ethics, Risks, and Use of AI Assistants 734

This work investigates the use of LLMs as expla- 735
nation interfaces for deterministic material plan- 736
ning systems, rather than as autonomous decision- 737
making agents. The proposed framework does 738
not alter MRP decisions, forecasts, or planning 739
logic, but generates natural-language explanations 740
grounded in an explicit, verifiable execution trace. 741

Potential risks arise if generated explanations 742
are interpreted as normative recommendations or 743
if trace information is incomplete or incorrectly 744
specified. To mitigate these risks, the proposed 745
system enforces strict schema normalization, deter- 746
ministic trace exposure, and evidence-constrained 747
generation, explicitly restricting explanations to 748
trace-derived facts. 749

AI assistants were used during the research pro- 750
cess for language refinement and editing, but all 751
technical design choices, formal definitions, exper- 752
imental setups, and interpretations were developed 753
and verified by the authors. The authors take full 754
responsibility for the content of the paper and its 755
claims. 756

References 757

- Abdulkadir Atalan. 2025. The ChatGPT application 758
on quality management: A comprehensive review. 759
Journal of Management Analytics, 12(2):229–259. 760
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 761
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 762
Plappert, Jerry Tworek, Jacob Hilton, and Reiichiro 763
Nakano. 2021. [Training verifiers to solve math 764](#)
[word problems](#). *Computing Research Repository*, 765
arXiv:2110.14168. 766
- Yilun Du, Shuang Li, Joshua B. Tenenbaum, and Anto- 767
nio Torralba. 2023. Improving factuality and reason- 768
ing in language models through multi-agent debate. 769
In *Proceedings of the 40th International Conference 770*
on Machine Learning. 771
- Faisal Mohamed Elbahri, Omar Ismael Al-Sanjary, 772
Musab A. M. Ali, Zakiya Ali Naif, Omar Ahmed 773
Ibrahim, and M. N. Mohammed. 2019. Difference 774
comparison of SAP, Oracle, and Microsoft solutions 775
based on cloud ERP systems: A review. In *2019 776*
IEEE 15th International Colloquium on Signal Pro- 777
cessing & Its Applications (CSPA). IEEE. 778

779	Bárbara Ferreira and João Reis. 2023. Artificial intelligence in supply chain management: A systematic literature review and guidelines for future research. In <i>Proceedings of the International Joint Conference on Industrial Engineering and Operations Management</i> , Cham. Springer.	832
780		833
781		834
782		
783		
784		
785	Jan C. Fransoo, Toni Waeffler, and John R. Wilson, editors. 2010. <i>Behavioral Operations in Planning and Scheduling</i> . Springer Science & Business Media, New York, NY.	
786		
787		
788		
789	David Gunning, Mark Stefik, Jaesik Choi, Timothy Miller, Simone Stumpf, and Guang-Zhong Yang. 2019. XAI—explainable artificial intelligence. <i>Science Robotics</i> , 4(37):eaay7120.	
790		
791		
792		
793	Dávid Gyulai, András Pfeiffer, and László Monostori. 2017. Robust production planning and control for multi-stage systems with flexible final assembly lines. <i>International Journal of Production Research</i> , 55(13):3657–3673.	
794		
795		
796		
797		
798	Wallace J. Hopp and Mark L. Spearman. 2011. <i>Factory Physics</i> . Waveland Press, Long Grove, IL.	
799		
800	F. Robert Jacobs, Richard B. Chase, and Nicholas J. Aquilano. 2011. <i>Manufacturing Planning and Control for Supply Chain Management</i> . McGraw-Hill, New York, NY.	
801		
802		
803		
804	Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4198–4205. Association for Computational Linguistics.	
805		
806		
807		
808		
809		
810	Scott M. Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In <i>Advances in Neural Information Processing Systems</i> .	
811		
812		
813	Vincent A. Mabert. 2007. The early road to material requirements planning. <i>Journal of Operations Management</i> , 25(2):346–356.	
814		
815		
816	Abhilash Puthanveetil Madathil, Xichun Luo, Qi Liu, Charles Walker, Rajeshkumar Madarkar, and Yi Qin. 2025. A review of explainable artificial intelligence in smart manufacturing. <i>International Journal of Production Research</i> , pages 1–44.	
817		
818		
819		
820		
821	Joseph A. Orlicky. 1974. <i>Material Requirements Planning: The New Way of Life in Production and Inventory Management</i> . McGraw-Hill, New York, NY.	
822		
823		
824	Afshin Oroojlooyjadid, MohammadReza Nazari, Lawrence Snyder, and Martin Takáč. 2022. A deep Q-network for the beer game: A deep reinforcement learning algorithm to solve inventory optimization problems. <i>Manufacturing & Service Operations Management</i> , 24(1):285–304.	
825		
826		
827		
828		
829		
830	Carol A. Ptak. 2011. <i>Orlicky’s Material Requirements Planning</i> , 3 edition. McGraw-Hill, New York, NY.	
831		
	Carol A. Ptak and Chad Smith. 2016. <i>Demand Driven Material Requirements Planning (DDMRP)</i> . Industrial Press, New York, NY.	832
		833
		834
	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?": Explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining</i> , pages 1135–1144.	835
		836
		837
		838
		839
		840
	David Salinas, Valentin Flunkert, and Jan Gasthaus. 2020. DeepAR: Probabilistic forecasting with autoregressive recurrent networks. <i>International Journal of Forecasting</i> , 36(3):1181–1191.	841
		842
		843
		844
	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 68539–68551.	845
		846
		847
		848
		849
		850
	Hartmut Stadler. 2005. Supply chain management and advanced planning: Basics, overview and challenges. <i>European Journal of Operational Research</i> , 163(3):575–588.	851
		852
		853
		854
	Daniel C. Steele. 1975. The nervous MRP system: How to do battle. <i>Production and Inventory Management</i> , 16(4):83–89.	855
		856
		857
	Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. 2023. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. In <i>Advances in Neural Information Processing Systems</i> , volume 36.	858
		859
		860
		861
		862
	Guodong Wang, Huabing Wang, Shanshan Lv, Xinxue Kang, and Qing’an Cui. 2025. Explainable artificial intelligence for manufacturing process modelling and optimisation: An integrated ANN and LIME approach. <i>International Journal of Production Research</i> , pages 1–22.	863
		864
		865
		866
		867
		868
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 24824–24837.	869
		870
		871
		872
		873
		874
	Qingyun Wu, Gagan Bansal, Jason Wei, Yi Zhang, Yiran Wu, Qingsong Wen, Andrew M. Dai, and Denny Zhou. 2024. AutoGen: Enabling next-gen LLM applications via multi-agent conversations. In <i>Proceedings of the First Conference on Language Modeling</i> .	875
		876
		877
		878
		879
	Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2024. Explainability for large language models: A survey. <i>ACM Transactions on Intelligent Systems and Technology</i> , 15(2):1–38.	880
		881
		882
		883
		884

885 A Formal Definitions and Notation

886 This appendix defines the core objects used in
887 EXPLAINABLEMRP to ensure clarity and repro-
888 ducibility.

889 **MRP Inputs.** For a material m , the input in-
890 stance is defined as $X_m = \{S_0, \mathcal{D}, \mathcal{P}, L, \mathcal{H}\}$,
891 where S_0 is the initial on-hand stock, $\mathcal{D} =$
892 $\{(t_i, d_i)\}$ is the set of dated demand events, $\mathcal{P} =$
893 $\{(t_j, p_j)\}$ is the set of dated supply (purchase or-
894 der) events, L is a fixed lead time, and \mathcal{H} is a dis-
895 crete planning horizon.

896 **Projected Inventory and Shortage.** Projected
897 inventory at time t is computed by deterministic
898 netting:

$$899 S(t) = S_0 + \sum_{t_j \leq t} p_j - \sum_{t_i \leq t} d_i. \quad (3)$$

900 The first shortage date is:

$$901 t^* = \min\{t \in \mathcal{H} : S(t) < 0\}. \quad (4)$$

902 **Deterministic MRP Trace.** A deterministic trace
903 for material m is a time-ordered record of netting
904 states:

$$905 \mathcal{T}_m = \{(t, S(t), D_t, P_t) \mid t \in \mathcal{H}\}, \quad (5)$$

906 where $D_t = \sum_{t_i \leq t} d_i$ is cumulative demand and
907 $P_t = \sum_{t_j \leq t} p_j$ is cumulative supply up to time t .
908 Each element of the trace corresponds to an explicit
909 instantiation of Eq. 3 at time t , and the shortage
910 condition is detected by Eq. 4.

911 **Evidence Extraction.** We define a deterministic
912 evidence extractor

$$913 f_{ev} : \mathcal{T}_m \rightarrow E_m, \quad (6)$$

914 where E_m is a finite set of trace-validated facts used
915 as admissible support for explanations. In our im-
916 plementation, evidence includes (non-exhaustively)
917 the shortage date, total demand/supply, usage spike
918 dates, and a recent inventory trend, all computed
919 deterministically from the trace (see Appendix B
920 for an example).

921 **Faithfulness and Trace Alignment Require-**
922 **ments.** Let y_m be a natural-language explana-
923 tion for material m , and let $I(y_m)$ be the set of
924 *informational atoms* extracted from y_m (e.g., dates,
925 quantities, and causal relations). We require the
926 faithfulness condition:

$$927 I(y_m) \subseteq E_m. \quad (7)$$

928 In addition, every causal statement in y_m must cor-
929 respond to a verifiable step in \mathcal{T}_m (*trace alignment*),
930 so that explanations reflect the deterministic netting
931 dynamics rather than post hoc rationalizations.

932 **Faithfulness Metrics.** We evaluate faithfulness
933 using Evidence Coverage (EC) and Hallucination
934 Rate (HR).

935 Let $\text{EvidenceUsed}(y_m) \subseteq E_m$ denote the sub-
936 set of evidence items referenced in y_m .

$$937 \text{EC}(y_m) = \frac{|\text{EvidenceUsed}(y_m)|}{|E_m|}. \quad (8)$$

938 Let $\text{Hallucinated}(y_m)$ be the set of statements in
939 y_m that cannot be grounded in either \mathcal{T}_m or E_m .

$$940 \text{HR}(y_m) = \frac{|\text{Hallucinated}(y_m)|}{|\text{Statements}(y_m)|}. \quad (9)$$

941 **Trace Alignment Score (TAS).** TAS is a human
942 evaluation score assessing whether causal claims
943 in y_m are consistent with the deterministic trace
944 \mathcal{T}_m and its implied netting dynamics. Evaluators
945 rate *faithfulness*, *specificity*, and *clarity* on a Likert
946 scale; we report averaged scores across evaluators
947 and materials in the main text.

948 B Example of a Deterministic MRP Trace

949 Table 4 shows an excerpt of the deterministic MRP
950 trace for *Material 1*. The trace records cumulative
951 demand D_t , cumulative applied PO P_t , and the re-
952 sulting projected inventory $S(t)$ over the planning
953 horizon. In our implementation, when the inven-
954 tory first becomes negative, the system registers
955 the shortage date and immediately applies all re-
956 maining open PO quantity in the same step (see the
957 *SimulationAgent* logic).

958 This example illustrates two key properties.
959 First, the shortage event is uniquely determined
960 by the cumulative effect of dated demand and sup-
961 ply events (Eq. 3–4). Second, trace visibility makes
962 the causal sequence auditable: planners can iden-
963 tify the exact demand increment that first violates
964 $S(t) \geq 0$ and verify how supply actions affect
965 subsequent inventory evolution.

966 C Deterministic Execution and Evidence 967 Construction

968 This appendix documents the deterministic compo-
969 nents underlying EXPLAINABLEMRP, including
970 the raw input data format, the MRP execution al-
971 gorithm, evidence extraction, and the prompt struc-
972 tures used for explanation generation.

Date	Demand _t	PO _t	S _{t-1}	S(t) before PO	S(t) after PO
2026-03-20	139.191	0	324.900	185.709	185.709
2026-03-21	14.143	0	185.709	171.566	171.566
2026-03-22	125.060	0	171.566	46.506	46.506
2026-03-23	139.203	276	46.506	-92.697	183.303
2026-03-24	14.143	0	183.303	169.160	169.160
2026-03-25	125.060	0	169.160	44.100	44.100
2026-03-26	139.263	0	44.100	-95.163	-95.163

Table 4: Excerpt of the deterministic netting trace for Material 1 around the first shortage event. Demand_t denotes the daily demand increment; PO_t denotes the PO quantity applied at that date. The first shortage is detected on 2026-03-23 when S(t) becomes negative before PO application, after which remaining PO is applied immediately.

973 C.1 Raw Dataset Format (Masked Industrial 974 Data)

975 Experiments are conducted on anonymized opera-
976 tional datasets derived from an industrial manufac-
977 turing environment. All identifiers and numerical
978 values are masked, while preserving the structural,
979 temporal, and causal relationships required for de-
980 terministic MRP execution.

981 The raw dataset consists of four enterprise
982 spreadsheets commonly found in ERP systems:

983 **Bill of Materials (BOM).** The BOM table speci-
984 fies component-to-product relationships and usage
985 coefficients. Columns include product identifiers,
986 component identifiers, usage quantities, and unit
987 information. Column names and units vary across
988 files and vendors.

989 **Production Plan (Plan).** The Plan table encodes
990 time-phased production demand. Each column may
991 correspond to a production date or batch, with het-
992 erogeneous naming conventions and mixed units.

993 **Inventory Records (Stock).** The Stock table pro-
994 vides initial on-hand inventory levels. It may in-
995 clude auxiliary fields such as warehouse locations,
996 inspection status, or free-text comments that are
997 not directly relevant to MRP computation.

998 **Open Purchase Orders (OpenPO).** The
999 OpenPO table lists outstanding supply orders.
1000 In addition to structured fields (quantity, due
1001 date), this table often contains unstructured
1002 natural-language comments (e.g., supplier delays
1003 or logistics issues).

1004 These spreadsheets are heterogeneous in schema,
1005 naming conventions, and units. Without explicit
1006 normalization, large language models frequently
1007 misinterpret column semantics, motivating the use
1008 of a deterministic schema-aware DataAgent.

C.2 Deterministic MRP Execution Algorithm 1009

The *SimulationAgent* implements a deterministic,
1010 rule-based MRP execution engine. This component
1011 performs exact netting and offset computation and
1012 does not invoke any language model. 1013

Inputs. 1014

- Initial on-hand stock S_0 1015
- Demand events $\mathcal{D} = \{(t_i, d_i)\}$ 1016
- Supply events $\mathcal{P} = \{(t_j, p_j)\}$ 1017
- Fixed lead time L 1018
- Discrete planning horizon $\mathcal{H} = \{t_1, \dots, t_T\}$ 1019

Output. A deterministic MRP trace 1020

$$\mathcal{T} = \{(t, S(t), D_t, P_t)\}_{t \in \mathcal{H}}. \quad 1021$$

This algorithm explicitly instantiates the inven-
1022 tory update rule (Eq. 1) at every time step. Given
1023 fixed inputs, the resulting trace \mathcal{T} is uniquely deter-
1024 mined. 1025

PO Application Convention. In our implemen-
1026 tation, when the inventory first becomes negative,
1027 all remaining open purchase order quantities sched-
1028 uled at the same time step are applied immediately.
1029 This convention is used consistently in Algorithm 1
1030 and the trace examples reported in Appendix B. 1031

C.3 Deterministic Evidence Extraction Function 1032

To ensure faithful explanation, explanations are re-
1034 stricted to a finite set of evidence deterministically
1035 extracted from the trace. 1036

We define an evidence extraction function 1037

$$f_{ev} : \mathcal{T} \rightarrow E, \quad 1038$$

where E denotes the admissible evidence set. 1039

Algorithm 1 Deterministic MRP Trace Computation

Require: Initial stock S_0 , demand events \mathcal{D} , supply events \mathcal{P} , planning horizon \mathcal{H} , lead time L

Ensure: Deterministic trace \mathcal{T}

```
1: cumulative_demand  $\leftarrow 0$ 
2: cumulative_supply  $\leftarrow 0$ 
3: inventory  $\leftarrow S_0$ 
4:  $\mathcal{T} \leftarrow \emptyset$ 
5: shortage_detected  $\leftarrow \mathbf{False}$ 
6: for each time  $t \in \mathcal{H}$  (chronological order) do
7:   for all  $(t_i, d_i) \in \mathcal{D}$  where  $t_i = t$  do
8:     cumulative_demand  $\leftarrow$  cumulative_demand  $+ d_i$ 
9:   end for
10:  for all  $(t_j, p_j) \in \mathcal{P}$  where  $t_j = t$  do
11:    cumulative_supply  $\leftarrow$  cumulative_supply  $+ p_j$ 
12:  end for
13:  inventory  $\leftarrow S_0 +$  cumulative_supply  $-$  cumulative_demand
14:  Append  $(t, \text{inventory}, \text{cumulative\_demand}, \text{cumulative\_supply})$  to  $\mathcal{T}$ 
15:  if inventory  $< 0$  and shortage_detected = False then
16:    shortage_date  $\leftarrow t$ 
17:    arrival_date  $\leftarrow t - L$ 
18:    shortage_detected  $\leftarrow \mathbf{True}$ 
19:  end if
20: end for
```

Evidence Items. From a trace \mathcal{T} , the extractor deterministically computes:

- Initial stock S_0
- Total cumulative demand D_T
- Total cumulative supply P_T
- First shortage date $t^* = \min\{t : S(t) < 0\}$
- Inventory trajectory $\{S(t)\}_{t \in \mathcal{H}}$
- Demand spike events (large single-period increases)
- Supply gaps (intervals with increasing demand but no incoming supply)

All evidence items are computed via fixed rules without statistical inference or language model involvement.

Role in Explanation. Let y be a generated explanation and $I(y)$ its extracted informational atoms. Faithfulness is enforced by the constraint:

$$I(y) \subseteq E.$$

C.4 Prompt Templates and Evidence Injection

This section documents the prompt structures used to generate explanations.

Raw-LLM Baseline Prompt. The baseline prompt provides raw spreadsheets and requests an explanation without trace access or constraints:

Explain why a shortage occurs for the given material based on the provided tables.

ExplainableMRP Prompt. The proposed framework uses a structured prompt consisting of:

1. Schema-normalized input tables
2. Deterministic trace-derived evidence serialized as JSON
3. An explicit instruction restricting explanations to the provided evidence

Evidence Constraint Instruction.

Generate an explanation using only the provided evidence. If a claim cannot be supported by the evidence, state that it is unavailable. Do not speculate or introduce external information.

Output Schema and Verification. The model must return a structured JSON object containing reasoning, cause, action, EvidenceUsed, and

1083 AlignmentIssues. Outputs are parsed deterministically to verify evidence usage and trace alignment.
1084
1085 If validation fails, the system falls back to a deterministic template-based explanation.
1086

1087 **Scope of Evaluation.** All quantitative and qualitative evaluations reported in this paper (RQ1–
1088 RQ3) are conducted on baseline MRP executions
1089 without hypothetical scenario modifications. While
1090 the system supports interactive scenario simulation
1091 (e.g., demand scaling or supply removal), these capabilities are treated as system features and are not
1092 evaluated in the current study.
1093
1094

1095 D Evaluation and Verification Details

1096 This appendix documents the parsing and validation logic used for faithfulness verification, as well
1097 as the human evaluation protocol and raw evaluation results. These details are provided to ensure
1098 transparency and to address potential concerns regarding subjectivity or heuristic bias.
1099
1100
1101

1102 D.1 Parsing Rules and Validation Logic

1103 To verify explanation faithfulness, all generated
1104 explanations are subjected to deterministic parsing
1105 and validation.

1106 **Informational Atom Extraction.** Given an explanation y , we extract a set of *informational atoms*
1107 $I(y)$ using rule-based pattern matching. Atoms correspond to explicitly stated: (i) dates (e.g., shortage
1108 dates or arrival dates), (ii) numerical quantities (e.g., inventory levels, demand or supply amounts),
1109 and (iii) causal relations linking demand or supply events to inventory changes.
1110
1111
1112
1113

1114 The extraction process relies on predefined lexical templates and regular expressions, rather than
1115 learned classifiers, to ensure deterministic behavior.
1116
1117

1118 **Unsupported Claim Detection.** An informational atom $a \in I(y)$ is considered *unsupported*
1119 if $a \notin E$, where E is the trace-derived evidence set. Unsupported atoms are flagged as
1120 AlignmentIssues. This criterion is purely rule-based and does not involve human judgment or
1121 language model inference.
1122
1123
1124

1125 **Evidence Usage Tracking.** Atoms successfully
1126 matched to evidence items are recorded in
1127 EvidenceUsed. This mapping enables automatic
1128 computation of EC and HR.

1129 **Fallback Handling.** If parsing fails due to malformed output or missing required fields, the system
1130 falls back to a deterministic template-based
1131 explanation that reports only the shortage date and
1132 aggregate demand and supply. This conservative
1133 fallback prevents unverified claims from entering
1134 the evaluation.
1135

1136 **Fallback Usage in Evaluation.** To assess
1137 whether fallback handling influenced the reported
1138 hallucination results, we manually inspected all explanations generated for RQ3 using the evaluation
1139 spreadsheet. We confirm that no explanations in
1140 the evaluated dataset were replaced by template-based
1141 fallbacks. All evaluated outputs contained
1142 trace-derived evidence items and material-specific
1143 numerical references. Accordingly, the reported
1144 zero HR reflects the effect of explicit evidence constraints and deterministic validation, rather than
1145 post-hoc output filtering or output substitution.
1146
1147

1148 D.2 Human Evaluation Protocol

1149 Human evaluation is conducted to assess whether
1150 trace visibility improves the perceived quality of
1151 explanations, particularly in terms of causal grounding.
1152

1153 **Material Selection.** The evaluated materials are
1154 selected to represent non-trivial planning cases in
1155 which deterministic netting produces meaningful
1156 causal structures. Specifically, materials are included if they (i) experience at least one shortage
1157 event within the planning horizon, (ii) involve multiple dated demand events, and (iii) have overlapping
1158 supply and demand periods. Materials with trivial configurations (e.g., a single demand event
1159 or immediate supply satisfaction) are excluded to
1160 avoid degenerate explanation cases.
1161
1162
1163

1164 **Evaluators.** Three domain-informed evaluators
1165 with experience in material planning and manufacturing operations participated in the study. Evaluators
1166 were not involved in system development.
1167

1168 **Evaluation Setup.** For each material, evaluators
1169 are presented with two explanations corresponding to identical MRP outcomes: one generated with
1170 trace visibility and one without. The presentation order is randomized for each evaluator and each
1171 material. Evaluators are not informed which explanation has access to the deterministic trace.
1172
1173
1174

1175 **Evaluation Interface.** The evaluation is implemented using a web-based survey tool. For each
1176

1177	explanation pair, evaluators answer the following	Randomization Protocol. To mitigate ordering	1219
1178	questions:	and presentation bias:	1220
1179	• Faithfulness: Is the explanation consistent	• The assignment of explanations to labels <i>A</i>	1221
1180	with deterministic MRP logic?	and <i>B</i> was randomized for each question.	1222
1181	• Specificity: Does the explanation identify con-	• The presentation order of explanation condi-	1223
1182	crete dates, quantities, and causal events?	tions (e.g., with-trace vs. no-trace) was ran-	1224
1183	• Clarity: Is the explanation easy to understand	domized across questions.	1225
1184	and well structured?	Participants were not informed about how the	1226
1185	Each dimension is rated on a three-point Likert	explanations were generated or about the internal	1227
1186	scale.	mechanisms of the system.	1228
1187	In addition, evaluators indicate an overall prefer-	Participant Instructions. Participants were in-	1229
1188	ence between the two explanations.	structed to:	1230
1189	Blinding and Randomization. All evaluations	• Read both explanations carefully.	1231
1190	are double-blinded with respect to explanation con-	• Evaluate which explanation better explains	1232
1191	dition. Randomization is applied independently per	the shortage situation based solely on the pro-	1233
1192	material to mitigate ordering effects.	vided summary and trace.	1234
1193	D.3 Human Evaluation Instructions	• Rely on intuitive and honest judgment; no	1235
1194	This subsection provides the full instructions and	prior expertise in MRP or supply-chain plan-	1236
1195	questionnaire used for the human evaluation of	ning was required.	1237
1196	explanation quality.	Evaluation Criteria. Each explanation was eval-	1238
1197	Purpose of the Study. This survey was con-	uated independently using a 3-point Likert scale	1239
1198	ducted for research purposes to evaluate the quality	according to the following criteria.	1240
1199	of natural-language explanations for material short-	Faithfulness. <i>Is the explanation grounded in the</i>	1241
1200	age situations in MRP. All responses were collected	<i>provided summary and trace?</i>	1242
1201	anonymously and used solely for research.	• 1: Largely inconsistent with the evidence.	1243
1202	Overview of Each Evaluation Task. In each	• 2: Partially grounded in the evidence.	1244
1203	question, participants were presented with the fol-	• 3: Accurately reflects key quantities and tem-	1245
1204	lowing information for a single material:	poral dynamics.	1246
1205	• A summary of the shortage outcome (Sum-	Specificity. <i>Does the explanation refer to con-</i>	1247
1206	mary),	<i>crete quantities, dates, and changes?</i>	1248
1207	• A partial execution trace representing the ac-	• 1: Mostly abstract or vague.	1249
1208	tual deterministic computation that led to the	• 2: Mentions some numerical details.	1250
1209	shortage (Trace),	• 3: Clearly specifies dates, quantities, and	1251
1210	• Two natural-language explanations (Explana-	changes.	1252
1211	tion A and Explanation B).	Clarity. <i>Is the explanation easy to understand</i>	1253
1212	The execution trace serves as the ground-truth	<i>for non-experts?</i>	1254
1213	evidence indicating how the shortage occurred. Par-	• 1: Difficult to understand.	1255
1214	ticipants were asked to evaluate how well each	• 2: Moderately clear.	1256
1215	explanation aligns with the provided trace. Both	• 3: Very clear and easy to understand.	1257
1216	explanations were generated based on the same		
1217	material and the same shortage outcome, but may		
1218	differ in wording and structure.		

Overall Preference. After rating both explanations, participants were asked to indicate their overall preference:

- Explanation A,
- Explanation B,
- No clear preference.

Ethical Considerations. No personally identifiable information was collected. Participants were informed that their responses would be anonymized and used only for academic research purposes.

D.4 Raw Human Evaluation Scores

Table 5 reports the complete raw human evaluation scores for all materials, evaluators, and experimental conditions. Scores are reported without aggregation to allow independent inspection of per-material and per-evaluator variation.

D.5 Statistical Analysis and Robustness Checks

Using the raw faithfulness scores reported in Table 5, we conduct a paired Wilcoxon signed-rank test across all material–evaluator pairs ($n = 24$). The test yields a statistic of $W = 59.5$ with a p -value of 0.62, indicating that the aggregate difference between the trace and no-trace conditions is not statistically significant.

This result does not imply the absence of systematic benefits from trace visibility. As discussed in the main text, the effect of trace access varies substantially across materials with different structural complexity, which is not fully captured by aggregate significance tests.

Paired t -test as a Robustness Check. To complement the non-parametric analysis, we additionally perform a paired t -test on the same material–evaluator pairs. Table 6 summarizes the results.

Across all three dimensions, the paired t -tests similarly indicate no statistically significant differences, with small effect sizes (Cohen’s $d_z \leq 0.17$). These results are consistent with the Wilcoxon test and reinforce the conclusion that trace visibility does not yield uniform improvements across all materials.

Interpretation. The statistical analyses indicate that the effect of trace visibility is not uniform across all materials. Aggregate significance tests therefore provide an incomplete characterization of

Mat.	Eval.	Condition	Faith.	Spec.	Clar.
1	E1	No-Trace	2	2	3
	E1	With-Trace	3	3	2
	E2	No-Trace	2	2	3
	E2	With-Trace	3	3	2
	E3	No-Trace	2	2	2
	E3	With-Trace	2	2	3
2	E1	No-Trace	2	2	3
	E1	With-Trace	2	2	3
	E2	No-Trace	2	2	3
	E2	With-Trace	2	3	3
	E3	No-Trace	2	2	2
	E3	With-Trace	2	2	2
3	E1	No-Trace	2	2	2
	E1	With-Trace	3	3	3
	E2	No-Trace	1	2	1
	E2	With-Trace	2	3	3
	E3	No-Trace	2	2	2
	E3	With-Trace	2	2	3
4	E1	No-Trace	2	2	2
	E1	With-Trace	2	2	3
	E2	No-Trace	2	2	2
	E2	With-Trace	2	3	3
	E3	No-Trace	2	2	3
	E3	With-Trace	2	2	2
5	E1	No-Trace	2	2	2
	E1	With-Trace	3	2	3
	E2	No-Trace	2	3	2
	E2	With-Trace	2	2	3
	E3	No-Trace	1	2	3
	E3	With-Trace	2	2	2
6	E1	No-Trace	3	3	3
	E1	With-Trace	2	2	3
	E2	No-Trace	3	3	3
	E2	With-Trace	2	2	2
	E3	No-Trace	3	2	3
	E3	With-Trace	3	2	3
7	E1	No-Trace	2	2	2
	E1	With-Trace	3	3	3
	E2	No-Trace	2	3	3
	E2	With-Trace	3	3	3
	E3	No-Trace	2	2	2
	E3	With-Trace	3	3	3
8	E1	No-Trace	3	3	3
	E1	With-Trace	2	2	2
	E2	No-Trace	3	3	3
	E2	With-Trace	2	2	2
	E3	No-Trace	2	3	3
	E3	With-Trace	1	2	2

Table 5: Complete raw human evaluation scores for all materials, evaluators, and conditions. Faith., Spec., and Clar. denote faithfulness, specificity, and clarity, respectively. All scores are reported on a five-point Likert scale.

its impact. As demonstrated by the qualitative anal-

Metric	Δ Mean	$t(23)$	p
Faithfulness	+0.083	0.49	0.63
Specificity	+0.083	0.44	0.66
Clarity	+0.167	0.81	0.43

Table 6: Paired t -test results comparing With-Trace and No-Trace conditions ($n = 24$).

ysis, trace-based explanations are especially valuable in materials with deeper causal structures and longer planning horizons, where causal verification is otherwise difficult. In contrast, for materials with relatively linear demand–supply dynamics, trace access yields limited additional benefit and may not translate into higher aggregate scores. These findings suggest that the primary contribution of trace visibility lies in supporting causal verification in complex planning scenarios, rather than in consistently improving average evaluation metrics.

D.6 Qualitative Examples

To complement quantitative analysis, we include representative qualitative examples illustrating when trace visibility provides clear explanatory benefits and when its impact is limited. All examples are reproduced verbatim from the evaluated explanations.

D.7 Qualitative Failure and Correction Example

We present a representative qualitative example to illustrate when trace visibility materially improves causal verification.

Material 3 (Delayed Replenishment). In the no-trace condition, the explanation incorrectly assumed an immediate inventory recovery following a purchase order, despite a non-zero lead time. With trace access, the explanation explicitly referenced the delayed arrival date recorded in the trace, correctly attributing the shortage to cumulative demand growth during the lead-time window.

This example highlights how trace visibility prevents superficially fluent but causally incorrect explanations by exposing temporally grounded constraints that are otherwise omitted.

E Frontend and System Implementation Details

This appendix documents the interactive frontend used in our experiments and demo. The frontend is implemented in Streamlit and orchestrates the full agent pipeline, including automatic file role

classification, schema normalization, deterministic MRP execution, and trace-grounded explanation generation. All screenshots shown in this appendix correspond to the actual system used in our study and are presented in Korean, reflecting the language of the underlying industrial deployment.

E.1 Frontend Overview

The frontend provides an end-to-end workflow: (i) multi-file spreadsheet upload, (ii) automatic assignment of files to MRP roles (BOM/Plan/Stock/OpenPO), (iii) in-browser table editing and lead-time specification, (iv) baseline MRP execution and optional what-if scenario simulation, and (v) inspection of deterministic traces and agent logs.

The interface is explicitly designed to expose intermediate computation artifacts (e.g., trace tables, shortage dates, and evidence summaries), rather than only final shortage endpoints typically shown in commercial MRP systems.

E.2 Agent Orchestration in the Frontend

After file upload, the frontend applies an *AutoClassifierAgent* to infer the role of each spreadsheet (BOM, Plan, Stock, OpenPO). Users may manually override the assignment before execution. The selected tables are then passed to a *DataAgent* for schema normalization (*unify*), producing canonical BOM, Stock, and OpenPO tables for downstream deterministic computation.

Auxiliary AutoClassifierAgent. In addition to the three core agents described in the main paper (*DataAgent*, *SimulationAgent*, and *ReasoningAgent*), the frontend includes an auxiliary *AutoClassifierAgent* that assists users by automatically inferring the roles of uploaded spreadsheets. This component operates exclusively at the data ingestion stage and does not participate in deterministic MRP execution, evidence extraction, or explanation generation. Accordingly, it is excluded from all experimental evaluations reported in this paper.

MRP execution is performed by a deterministic runner (*MRPRunner*), which produces per-material summaries and full simulation traces. The frontend supports two execution modes:

- **Baseline:** execute deterministic MRP on the normalized inputs.
- **Scenario:** apply user-specified perturbations

1393 (pre-defined or free-form) and execute deter-
1394 ministic MRP on the modified inputs.

1395 **E.3 Scenario Specification and Application**

1396 The frontend supports both button-based and free-
1397 form natural language scenarios. Pre-defined sce-
1398 narios deterministically modify one of the follow-
1399 ing: (i) BOM usage (demand scaling), (ii) per-
1400 material lead times (lead-time shifts), or (iii) open
1401 purchase order quantities (supply scaling).

1402 For free-form scenarios, a language model is
1403 used only to parse the user’s natural-language
1404 instruction into a structured JSON specifica-
1405 tion with fields `material`, `leadtime_delta`,
1406 `demand_multiplier`, and `po_multiplier`. The
1407 resulting JSON is then applied deterministically
1408 to the relevant tables. Thus, language model us-
1409 age is strictly limited to scenario parsing, while all
1410 numerical MRP execution remains deterministic.

1411 **E.4 Trace and Log Exposure**

1412 To support explainability, the frontend exposes:

- 1413 • Per-material key dates (shortage date, required
1414 arrival date, suggested order),
- 1415 • A time-indexed simulation trace (inventory,
1416 cumulative demand, cumulative PO usage),
- 1417 • Demand spike markers derived from the trace,
- 1418 • Full agent logs, including `AutoClassifierA-`
1419 `gent`, `DataAgent`, `SimulationAgent`, and `Rea-`
1420 `soningAgent` logs, shown in expandable pan-
1421 els.

1422 Figure 2 illustrates the EXPLAINABLEMRP
1423 frontend, which exposes both intermediate agent
1424 reasoning traces and the final trace-grounded ex-
1425 planation. By making input interpretation, schema
1426 normalization, deterministic MRP execution, and
1427 shortage reasoning explicitly visible, the interface
1428 enables users to inspect and verify every step from
1429 raw data ingestion to explanation generation.

Agent Reasoning Trace and Deterministic Execution Logs

AutoClassifierAgent Logs (Automatic File Role Classification)

- File plan.xlsx classified based on column schema
- File openpo_explainable.xlsx classified with highest Open PO score
- File bom.xlsx assigned as BOM (score = 9)
- File stock.xlsx assigned as Stock (score = 6)

! (Shows automatic role inference for BOM / Plan / Stock / Open PO)

DataAgent Logs (Schema Normalization)

- Normalization and schema unification started
- Column renaming applied (e.g., Product → Material, Usage → Qty)
- Missing comment fields automatically inserted
- Canonical schema generated for deterministic execution

SimulationAgent Summary (Per-Material Outcome)

- Deterministic MRP executed for each material
- Shortage dates and key metrics computed
- Causally relevant events identified from trace

SimulationAgent Raw Logs (Detailed MRP Execution)

- Step-by-step inventory updates
- Explicit netting-and-offset computation
- Full deterministic trace available

ReasoningAgent Logs (LLM Explanation Generation)

- Explanation generated under evidence constraints
- Evidence usage verified
- Alignment issues reported when detected

Detailed View

Shortage Reasoning

Current inventory levels are critically low relative to total demand. There are no remaining open purchase orders to replenish stock, and inbound supply has been delayed due to a supplier shipment issue. As a result, cumulative demand cannot be satisfied, leading to a shortage.

Reasoning Steps (Trace-Aligned)

- The current inventory level is extremely low, indicating a shortage state.
- No additional inventory is secured because remaining purchase order quantities are zero.
- Inventory replenishment is delayed due to a supplier shipment delay.
- Cumulative purchase order usage is insufficient to meet total demand.
- To resolve this issue, urgent coordination with the supplier and additional ordering are required.

Key Dates

- Shortage Date: 2025-12-31
- Required Arrival Date: 2025-11-01
- Suggested Order Date: 2025-10-25
- Lead Time: 60 days
- Remaining Open PO: 0
- Supplier Comment: Supplier A shipment delay.

Recommended Action

Based on the supplier delay, immediate coordination with Supplier A is required to expedite shipment schedules and place additional purchase orders. Alternative suppliers should also be considered to mitigate supply risk.

The chart displays the MRP trace over time. The left y-axis represents Stock Level (0 to -3000), and the right y-axis represents Cumulative Values (0 to 3000). The x-axis shows dates from Jan 2025 to Nov 2027. The Stock Level (blue line) starts at 0 and drops to -3000 by Nov 2027. Cumulative Demand (orange line) rises to 3000. Cumulative PO Used (green line) rises to 2700. Demand Spike (red line) is at 0. Shortage Date (purple line) is at 2025-12-31.

Figure 2: EXPLAINABLEMRP frontend exposing agent reasoning traces and trace-grounded explanations (Korean-language interface). The interface displays logs from the AutoClassifierAgent, DataAgent, SimulationAgent, and ReasoningAgent to make input classification, schema normalization, and deterministic MRP execution explicit and inspectable. The resulting explanation combines natural-language reasoning with numerical evidence derived directly from the deterministic trace, including shortage timing, key dates, recommended actions, and inventory-demand trajectories. The shortage event corresponds exactly to the time step at which projected inventory becomes negative in the MRP trace.