# NetSuite N/LLM: How to Auto-Classify Refund Reasons

By houseblend.io · Published December 29, 2025 · 32 min read



## Executive Summary

Modern enterprises face an urgent need to **standardize and analyze refund reasons** across sales channels. Unstructured or inconsistent "refund reason" data – often free-text entered by customers or support agents – hampers analytics, revenue recovery, and fraud detection. Historically, organizations relied on manual coding or fixed lists of reason codes, which fail to capture nuanced explanations or scale with data volume. In contrast, **Oracle NetSuite's new N/LLM module** (a SuiteScript interface to OCI Generative AI) enables *auto-classification* of refund reasons at scale. By leveraging large language models (LLMs) and semantic embeddings within NetSuite, companies can automatically assign free-text refund comments to standardized categories. This report investigates how N/LLM can be applied to refund-reason data to enforce consistency, improve insights, and reduce manual workload. We review the business context of returns, examine technical approaches (traditional vs. AI-based), and survey N/LLM capabilities (text generation, embeddings, RAG) relevant to classification. Data from industry sources highlight the stakes (e.g. U.S. e-commerce return volumes and common return causes (Source: www.glencoyne.com) (Source: corp.narvar.com). We discuss a prototype pipeline: extracting refund reason text via SuiteQL/saved searches, generating embeddings or using LLM prompts, and auto-labeling each entry. Case examples illustrate benefits and challenges. We close with implications for governance, accuracy, and future NetSuite AI directions.

## Introduction and Background

Refunds and returns are a perennial cost center in retail and e-commerce. In 2024, U.S. retail returns were estimated at **$890 billion** (Source: www.glencoyne.com), a figure which includes lost revenue, processing costs, and inventory write-offs. Typical online return rates range from **15–30%** of orders (as high as 40% in apparel) (Source: www.glencoyne.com), meaning **nearly one in five purchases** may eventually be returned. Common reasons span "item didn't fit," "not as described," "damaged," "arrived late," or even "changed my mind" (Source: corp.narvar.com) (Source: corp.narvar.com). For example, a 2021 consumer survey found **42%** of returns were due to size or fit issues (Source: corp.narvar.com), and research shows apparel/footwear returns are dominated by fit problems (Source: link.springer.com). Any refund process typically records a *reason code* or *comment* from the customer. However, many systems use only a fixed list of generic codes (e.g. "Wrong Size", "Defective", "Other"), forcing employees or customers to shoehorn nuanced situations into broad buckets (Source: link.springer.com).

This free-form data problem undermines analytics and operations. Without standardization, finance and support teams cannot easily aggregate or trend "why" refunds occur. As one research broker noted, e-commerce platforms often "resort to a list of predefined causes" but these generic lists "do not capture the intrinsic differences" among products and reasons (Source: link.springer.com). In practice, customers may write detailed explanations in support tickets or notes, but such narrative is rarely leveraged systematically.

**Oracle NetSuite** is a leading cloud ERP system used by over 40,000 organizations (Source: houseblend.io) (Source: houseblend.io). It unifies financials, inventory, CRM, and more in one data platform – the so-called *"Suiteness"* (Source: houseblend.io). This centralized data (orders, returns, customer histories) provides rich potential for analytics. NetSuite historically lacked built-in AI, but in 2024–2025 it introduced the **N/LLM SuiteScript module** for generative AI (Source: docs.oracle.com) (Source: houseblend.io). N/LLM lets developers call large language models (LLMs) from within NetSuite scripts, enabling tasks like text summarization, question-answering, and – crucially – classification or tagging of records. This integration allows companies to apply AI to ERP data *in situ*, rather than exporting to an external system (Source: houseblend.io). The potential is to automatically analyze free-text fields (such as refund comments) against the company's own data and taxonomy.

This report explores **"standardizing refund reasons with NetSuite N/LLM: auto-classification at scale."** We first review how refund reasons are handled today and why standardization matters. Then we examine how AI-based classification – and specifically NetSuite's N/LLM – can solve these problems. We detail technical approaches (embeddings, prompts, retrieval-augmented generation) and consider performance and governance issues. Case examples illustrate the concept. Finally, we discuss implications for operations and future enhancements. All claims are supported by industry data and technical documentation.

# Refund Reasons in Business Context

## The Cost and Prevalence of Returns

Returns and refunds represent significant operational and financial challenges. Retail returns not only reduce revenue but also incur logistics, labor, and inventory costs (Source: www.glencoyne.com). For example, **Retail returns reached ~$890B in 2024** in the U.S. alone (Source: www.glencoyne.com), a figure that underscores the scale. Many returned items cannot be resold at full price; one study found only ~50% of returned apparel can be liquidated at full cost (Source: www.glencoyne.com). In aggregate, the "true" cost of returns (shipping, labor, restocking fees, discounts, etc.) often **exceeds the refund amount itself** (Source: www.glencoyne.com), quietly eroding profits.

Typical return rates vary by industry. A research survey observed that **online apparel returns average ~20%** of purchases (vs ~9% offline) (Source: link.springer.com). Similarly, consumer electronics return rates might hover in the high single digits. Despite liberal "return policies" boosting sales, industry experts note that **10–15% of e-commerce sales** ultimately get returned (Source: topmostlabs.com). High return rates are especially common in fashion (due to sizing/fits) and in marketplaces where quality varies. With billions in play, even a few percentage points of "leakage" through avoidable returns are material.

## Common Refund Reasons

Understanding *why* returns happen is crucial for reducing them and protecting revenue. Studies have consistently identified a handful of leading return reasons. For instance:

- **Size/Fit Issues:** Clothing that "doesn't fit" is the top cause in many surveys. Narvar reports 42% of customers cited fit/size problems for their last return (Source: corp.narvar.com).Similarly, academic studies note footwear and apparel see much more size-related returns than defects (Source: link.springer.com).

- **Product Not as Described:** Misleading descriptions or photos lead to dissatisfaction. Consumers often return items if color, quality or features don't match the online listing (Source: corp.narvar.com).

- **Damaged/Defective Goods:** Physical damage during shipping or manufacturing defects prompt returns (e.g. "screen is cracked," "item not working"). Research indicates around 5% of electronics returns come from defects (Source: link.springer.com).

- **Logistics Errors:** Late delivery, wrong items, or incomplete shipments cause returns. Narvar highlights "arrived late" and "wrong item" among top reasons (Source: corp.narvar.com).

- **Value/Preference:** Some returns occur because a customer simply changes mind ("I just don't like it") or feels the item is not worth its price (Source: corp.narvar.com).

Each of these yield characteristic keywords: "too small," "wrong color," "broke," "never arrived," "just dislike," etc. In practice, e-commerce sites often ask customers to select a return code from a drop-down (e.g. "Item Damaged," "Incorrect Item," "Changed Mind"). But these canned lists are coarse. A recent research paper notes that **return reason lists are generic and applied across all categories**, ignoring product-specific nuances (Source: link.springer.com). For example, "size issues" is critical for shoes/clothes but irrelevant for electronics. The generic approach means companies lose detail, and raw feedback in free-text is lost.

Data on frequency by reason is sparse, but trade surveys give clues. For example, one report found 42% of respondents cited sizing issues (Source: corp.narvar.com), and Narvar's 2025 article lists seven "common reasons" including fit, description mismatch, damage, lateness, wrong item, and subjective dissatisfaction (Source: corp.narvar.com) (Source: corp.narvar.com). These sources indicate that returns cluster into a manageable number of categories, making classification feasible – *if* the data can be standardized.

## Challenges in Handling Refund Reason Data

Despite containing valuable insights, refund-reason fields are notoriously messy. Common issues include:

- **Free-Text Variation:** Customers or support reps often enter narrative comments ("It was too short," "Product smelled like chemicals"). Natural language is highly variable.
- **Inconsistent Codes:** Dropdown codes differ by retailer and are sometimes mis-used by staff or customers (e.g. selecting "defective" to avoid return shipping costs) (Source: patents.google.com).
- **Lack of Context:** A code like "Defect" doesn't specify if the item broke or simply didn't work as expected.
- **Data Silos:** Returns may be initiated on marketplaces (e.g. Amazon) then logged in ERP manually, losing original reason phrases.
- **Volume and Scale:** High-volume sellers (thousands of returns per day) cannot manually review each reason.

These problems frustrate analytics. For instance, finance teams cannot easily query "how many returns were due to quality issues this quarter" if "quality issues" is an inferred concept not present in structured fields. As Palma.ai observes generally for AI projects, "70–80% of AI projects never make it to production" due to data issues (Source: palma.ai). In the ERP context, poorly standardized fields are a common Achilles' heel. Lack of governance ("shadow IT" collecting reasons in untracked ways) only worsens the mess (Source: palma.ai).

Thus, there is strong motivation to **improve consistency** of refund-reason data. Standardized categories enable clear KPIs ("X% of returns due to defects"), automated decision rules (e.g. fraud detection for unusual reason patterns), and richer insights (e.g. root-cause analysis across products). Achieving this requires either forcing structured input or *automatically classifying* unstructured entries. The former (forcing) often alienates customers; the latter (AI-based classification) has become practical thanks to new tools. The next sections examine how AI, and specifically NetSuite's N/LLM, can automate this classification task at scale.

# Traditional vs. AI-Based Classification

Before exploring NetSuite's solution, we survey existing approaches to refund-reason standardization.

## Manual and Rule-Based Approaches

Traditionally, companies might use one or more of these methods:

- **Fixed Reason Codes:** As noted, many ERP systems allow defining a finite set of "return reasons" (like "Damaged," "Ordered by Mistake," etc.) for users to pick. While easy to implement, this rigid taxonomy often leads to user frustration and mislabeling{[53†L1-L4]}. It cannot account for nuanced or new reasons without reconfiguring the list.
- **Guided Forms:** Some e-commerce platforms add guided forms (e.g. checkboxes, radiobuttons) for returns, funneling customers through predetermined options. This can improve consistency but still misses nuance unless expanded ad infinitum. It also relies on the customer's understanding.
- **Post-hoc Manual Coding:** In some organizations, analysts periodically review a sample of free-text reasons and manually tag them into categories. This is labor-intensive and becomes impractical if returns volume is high.

- **Rule-based Parsing:** Simple natural language processing (NLP) scripts can flag keywords (e.g. if "broken" in comment, classify as "Defect"). While automatable, these brittle systems miss synonyms and context, and require constant maintenance. For example, "not working" vs "broken" vs "doesn't turn on" may all mean defective, but a rule parser needs each variant listed.

These methods often yield **low scalability or poor accuracy**. Manual categorization is costly and slow, while static rules cannot handle the diversity of customer language. The heterogeneity of text literally grows with more customers and more catalogs. As RecordPoint observes in the broader context, "manual classification scales poorly when organizations generate terabytes of new data daily" (Source: www.recordpoint.com). In practice, companies that rely on basic methods often find their return analysis reports are incomplete or inconsistent.

## Machine Learning and NLP Methods

In recent years, data-driven approaches have become viable:

- **Traditional ML Classification:** One could collect historical returns labeled by reason, then train a supervised classifier (e.g. logistic regression, SVM, or even a neural net) on TF-IDF features. These models can parse subtle variations better than keyword rules. Academic research on complaint classification or review sentiment provides analogies (Source: ashwin-ks.github.io) (Source: houseblend.io). However, such models need substantial labeled data (which may not exist initially) and can struggle with evolving language (require retraining).

- **Topic Modeling / Unsupervised Clustering:** Techniques like LDA or k-means on word embeddings can automatically find themes in return comments (as in the Springer study of return reasons from reviews (Source: link.springer.com). These methods can *suggest* categories by clustering similar phrases (e.g. "too small" and "runs tight" cluster together). While useful for discovery, they typically require human interpretation of clusters and may not align perfectly with business categories.

- **Embedding + Similarity:** A newer approach is to leverage pretrained embeddings: map each free-text reason into a high-dimensional vector (via models like BERT or Cohere) and then classify by nearest neighbor or vector clustering. This combines unsupervised semantic encoding with the ability to use labeled "seed" examples of each class. NetSuite's N/LLM offers an *embedding API* specifically for this purpose (Source: docs.oracle.com).

Machine learning has clear advantages: it can adapt to nuances (e.g. synonyms, conversational tone) and improve with more data. Indeed, specialized large-language models have been shown to **outperform traditional NLP** on many categorization tasks. For instance, LLMs trained on product data can extract attributes more accurately than older models (Source: houseblend.io). In-house studies (e.g. at large retailers) have demonstrated that ML-based tagging of unstructured fields yields far better recall/precision than static lists. The Amazon patent we examined explicitly uses ML to "estimate root causes" of returns from customer-entered data (Source: patents.google.com) – a recognition that free-text analysis is necessary to supplement flawed codes.

However, ML models traditionally require data science effort and infrastructure. They may also struggle to handle unknown categories gracefully or to "explain" their reasoning. This is where LLMs and RAG (retrieval-augmented generation) come in, promising a more flexible yet enterprise-grounded approach.

# NetSuite's Platform and the N/LLM Module

## NetSuite's Data Platform ("Suiteness")

Oracle markets NetSuite's unified data model as "**Suiteness**," i.e. having all ERP modules on one cloud platform (Source: houseblend.io). Because orders, returns, items, and customer interactions all live in the same environment, NetSuite accounts can easily aggregate and correlate data. This vertical integration is a boon for AI: as NetSuite's development leaders note, it means **"more of the data across your whole business"** is available for AI analysis (Source: houseblend.io). For example, BirdRock Home (a large home goods retailer on NetSuite) processes thousands of orders daily; their entire transaction history and product metadata reside in NetSuite (Source: houseblend.io). This wealth of labeled data (e.g. historical refund records) can fuel training and validation for classification solutions.

Pragmatically, NetSuite provides SuiteScript APIs to extract data. Developers can run **SuiteQL** (SQL-like queries) or Saved Searches to retrieve record fields. These can feed into AI pipelines either by passing them to an LLM prompt or using them as "documents" in RAG. Houseblend's analysis emphasizes that the key to N/LLM is efficiently surfacing *relevant* data (e.g. recent returns, item attributes) and supplying it as context (Source: houseblend.io). Because SuiteQL queries can join orders, items, and returns easily, an LLM prompt could theoretically include, say, "Customer X returned item Y (electronics) on date Z; follow-ups?".

NetSuite also offers the **Analytics Warehouse (NAW)**: a read-only data lake optimized for reporting and classical machine learning. NAW can integrate with Oracle's ML tools for large-scale analysis (like forecasting). But NAW is separate from the transactional NetSuite and is synced daily. In contrast, N/LLM acts *inside* NetSuite in real time. So for on-demand tasks like classifying a batch of refunds during nightly processing, N/LLM (via SuiteScript) is ideal. It uses the same data model and respects NetSuite governance without complex ETL. As Houseblend notes, N/LLM essentially lets developers call LLMs "as part of any script" with built-in support for prompts, embeddings, and source documents .

## The N/LLM Module: Capabilities Summary

The NetSuite **N/LLM** SuiteScript 2.1 module provides a set of functions to interact with OCI GenAI. Key capabilities include (Source: docs.oracle.com) (Source: docs.oracle.com):

- **Text Generation (LLM.chat/generateText):** Send a prompt to generate or complete text. Useful for tasks like creating summaries, writing descriptions, or answering questions based on prompt context (Source: houseblend.io). (For classification, one could craft prompts like "Categorize this reason:" and let the LLM respond with a label.)

- **Contextual Documents (RAG):** Using `llm.createDocument()`, developers can package strings (e.g. saved search results) as "documents" that the LLM will use for context. `generateText(options)` accepts an array of such documents. The model will incorporate that data into its answer and can even cite which document helped (Source: docs.oracle.com) . This enables *retrieval-augmented generation*: grounding the LLM's response in actual ERP records. For example, you could attach the returned item's category, vendor notes, and past return history as RAG context when classifying its refund reason.

- **Embeddings (llm.embed):** Convert arbitrary text (a refund comment, for instance) into a fixed-length vector (1536 dimensions with Cohere Embed v4.0 (Source: docs.oracle.com). The SuiteScript API can return these embeddings, which can then be compared via cosine similarity or input into downstream clustering/classification algorithms. Oracle explicitly notes embeddings are "useful for … text classification or text clustering" (Source: docs.oracle.com). In practice, one might embed each refund comment and compare it to prototype embeddings for each category, or cluster all refunds to discover emergent groups.

- **Prompt Management:** N/LLM supports a Prompt Studio where admins can define reusable prompts and parameters. Scripts can call these via `llm.evaluatePrompt` . This helps standardize the prompt format across the company. (For example, a "Refund Classifier" prompt could be set up once and invoked with each new comment.)

- **Streaming:** For interactive use cases, `generateTextStreamed` delivers model output token-by-token. This is less relevant for batch classification but could power real-time chatbots around refund inquiries.

The table below summarizes these capabilities:

| N/LLM CAPABILITY | DESCRIPTION | USE CASE | REFERENCE |
|---|---|---|---|
| **Text Generation (chat)** | Calls an LLM with a prompt; returns generated text (e.g. completion, answer). Supports parameters (model, etc.). | Summarizing notes, auto-drafting emails or reports, or answering Q&A from record data. | (Source: houseblend.io) |
| **Contextual Docs (RAG)** | Attach NetSuite data as "documents" for LLM context; the model uses them in output and can cite sources. | Grounded responses: e.g. answer "Why was this item returned?" using attached invoice, item details. | (Source: docs.oracle.com) (Source: houseblend.io) |
| **Embeddings** | Generates numeric vectors for input text for use in semantic search, clustering, or classification. | Cluster similar reasons, find nearest category embedding, or federate with ML classifier. | (Source: docs.oracle.com) |
| **Prompt Management** | Define reusable prompts with variables via Prompt Studio; scripts fill in values. | Consistent classification prompts or surveys for all returns. | (SuiteScript feature) |

In sum, N/LLM is a flexible toolkit. We will see how it can be assembled for refund-reason classification specifically.

# Auto-Classification of Refund Reasons Using N/LLM

In this section we describe how one would employ N/LLM tools to standardize refund reasons at scale. We break it down into steps: data collection, defining categories, model invocation, and evaluation.

## Data Collection in NetSuite

First, one gathers the relevant refund data from NetSuite. This typically involves:

- **Identifying the Records:** NetSuite has *Return Authorization* (RMA) records which may hold return reason information. If the company uses RMAs properly, each customer return has an RMA transaction (a non-posting record) tracking expected returns (Source: docs.oracle.com). That RMA can be linked to a Cash Sale or Invoice via `createdfrom` (Source: docs.oracle.com). In addition, a *Customer Refund* or *Credit Memo* (actual crediting or refunding transaction) may contain notes. In practice, a custom field might be used to record the textual reason selected by staff or customers. Even if not, support tickets or attached comments may hold free-text reasons.

- **SuiteQL / Saved Searches:** We can use SuiteScript to run a saved search or SuiteQL query on the return or refund records. For example, a SuiteQL query might pull the fields: return ID, returned item, date, refund amount, and *return reason note*. If reasons exist as a free-text field (`custbody_return_reason`) or there is an audit trail of message, we retrieve that text. Alternatively, one could export customer messages from email/chat integrated with NetSuite.

- **Data Format:** The result is effectively a table where each row is a return event with a "reason_text" field. This text might be one sentence from the customer, or a concatenation of notes. Some preprocessing (e.g. removing PII) may be needed.

According to NetSuite's governance model, these script queries run on the server (limited by governance units). Houseblend notes that efficient pipelining is crucial: one should carefully page or filter results for large datasets (Source: houseblend.io). In practice, a scheduled Suitelet script can collect, say, all unprocessed return comments and feed them to the classification step.

## Defining Categories (Taxonomy)

Parallel to data collection, the business must define the **target categories** for refund reasons. These could be a fixed list such as: *Defective; Sizing Issue; Wrong Item; Late Delivery; Not as Described; Other*. The categories should align with company objectives. For standardized reporting, one might group related causes (e.g. combining "wrong color" and "different style" under "Not as Described").

Each category can be represented in the model either via examples or descriptions. For instance, one might prepare a short description or a few sample phrases per category to use with the LLM. These will serve as a reference for classification. For example:

- *Defective*: "Item is broken, not working, damaged in shipping."
- *Sizing Issue*: "Too small, can't fit, size runs small, sleeves too short."
- *Wrong Item*: "I received a different product than ordered, wrong item sent."
- ... and so on.

This step injects domain knowledge into the process. The more representative the examples, the better the model can match them.

## Classification Methods

There are two main N/LLM-based approaches to assign a category to each reason text:

### 1. Embedding-Based Classification

Under this approach, we use the `llm.embed` API to convert text to vectors, then classify by similarity or clustering:

- **Generate Vectors:** Call `llm.embed({text: reason_text})` for each refund reason. This yields a 1536-dimensional embedding (for Cohere v4) for each text (Source: docs.oracle.com).

- **Reference Embeddings:** Also compute embeddings for category descriptions or prototypical examples prepared earlier. For example, embed the phrase "Size issue – runs small" for the *Sizing* category, etc.

- **Similarity Scoring:** For each reason vector, compute cosine similarity to each category prototype vector. Assign the category with highest similarity. This is a form of zero-shot classification. Because embeddings capture semantic meaning, "too hot, too small" will be near other *Sizing* examples. Oracle documentation explicitly mentions that embeddings can be used for "text classification or text clustering" (Source: [docs.oracle.com](docs.oracle.com)), which is exactly this.

- **Clustering (Optional):** Alternatively, cluster the embeddings of all reason texts (k-means or hierarchical clustering). Then a human can label clusters. Modern LLM embedding spaces often produce clearly grouped clusters (e.g. all "broken/damaged" texts cluster together). This is unsupervised but can suggest categories.

**Advantages:** Embeddings approach is efficient once the vectors are computed. We can reuse the category prototypes, and new items are classified by a quick cosine lookup. It avoids having to format complex prompts for each item. It also easily scales to large lists of reasons.

**Challenges:** It assumes pretrained embeddings capture the nuance of refund contexts. Rare or subtle categories might blur. Some disambiguation (e.g. "not working – user error vs defect") could be misclassified. One may need to tune thresholds or combine with heuristics. Also, embeddings do not "explain" classification decisions.

## 2. Generative Prompt-Based Classification

Here we formulate the task as a prompt to the LLM via `llm.generateText` (or `llm.chat`). For each reason, we ask the model explicitly to choose a category among a list. For example, a prompt might be:

```
System: You are an expert in analyzing retail return reasons.
Given the customer's reason for return and a list of categories, output the most appropriate category.
Categories: Defective, Sizing Issue, Wrong Item, Late Delivery, Not as Described, Other.

Customer Reason: "The shirt is way too small and sleeves are too short."
Category:
```

The LLM (e.g. Cohere, etc.) will generate an answer like "Sizing Issue".

We can use **Prompt Studio** to store this template and then call it with each reason via `llm.evaluatePrompt` (in which "reason" is a variable). This ensures consistency in formatting. *Retrieval*\*: We could also add RAG documents: e.g. attach a product description or previous customer comment, so the LLM knows the context of the item.

**Advantages:** The LLM has world knowledge and can handle nuance or ignore irrelevant text. For example, it might understand that "changed my mind" falls under "Other" or "Preference". It can parse complex sentences and even multiple sentences. No training is needed beyond prompt design.

**Challenges:** If done naïvely, this uses an LLM call per record, which could be slower and costlier than simple vector math. However, NetSuite's N/LLM may allow batch prompts (depending on the API usage). Also, we must ensure the prompt is designed to not hallucinate. We limit output by controlling options (maybe ask for exact string from list).

In NetSuite code, a SuiteScript could loop through each reason, call the prompt, and write back a custom field "Assigned Category = [result]".

Notably, Oracle's docs caution that generative responses require validation (Source: [docs.oracle.com](docs.oracle.com)). Thus, one should spot-check or use confidence indicators. Alternatively, one can use a two-tier system: run an initial pass with the LLM, then verify low-confidence cases manually, gradually feeding corrected labels back into a supervised model if needed.

## Example Workflow

Putting it together, a possible pipeline is:

1. **Extract refund data:** Use a saved search or SuiteQL to pull all recent return authorizations and their reason notes.

2. **Data cleansing:** Optionally normalize text (lowercase, strip punctuation). Remove entries with empty reasons.

3. **Define categories:** Finalize a list of 5–8 target labels that cover the business needs. Store any example phrases.

4. **Embeddings prep:** Call `llm.embed` on category keywords/phrases to get reference vectors.

5. **Classify each entry (embeddings):** For each reason text, call `llm.embed`, then compute cosine similarity to each category vector, assign highest-scoring category. Alternatively, group embeddings and label clusters for bulk classification.

   *OR*

6. **Classify each entry (prompt):** For each reason text, call `llm.generateText` or `llm.evaluatePrompt` with a prompt listing the categories and asking for the best match. Record the output category. (Include context if needed.)

7. **Store results:** Write the assigned category back to a custom field on the return authorization (or into a new classification table). This can then drive reporting.

8. **Validation:** Periodically check a random sample of classifications manually. If many are wrong, refine prompt phrasing or category definitions and rerun. Because NetSuite tracks script usage, one can monitor how many API calls are used and tune for efficiency.

An additional refinement could be using **LLM embeddings to retrain a lightweight model**. For example, once N/LLM has labeled thousands of reasons, the business could train a simpler in-ERP model (outside core scope) to classify reasons even faster. But the immediate goal is to leverage NetSuite's built-in AI.

# Data Analysis and Evidence

## Classification Performance Expectations

The key question is: *How accurate and actionable is N/LLM classification?* Since NetSuite's AI uses state-of-the-art models, we expect high semantic understanding. Industry evidence in related tasks is encouraging:

- **LLM Superiority:** As Houseblend noted, specialized LLMs outperformed traditional models on extracting product attributes (Source: houseblend.io). We similarly expect an LLM to correctly interpret a variety of return reason phrasings. For example, "too tight in waist" and "run small" both clearly hint at a sizing issue, which an LLM embedding should cluster together.

- **Human-Like Understanding:** N/LLM can use context. If we attach the item's category as context (e.g. electronics vs clothing), the LLM will apply relevant logic (electronics rarely have "fit" issues). This RAG grounding can reduce mistakes. In enterprise deployments, RAG has been shown to "ensure answers rely on factual context" 【6†L28-L34】. (Source: houseblend.io)ata from pilots would be ideal, but at time of writing such case studies are not publicly available. We can estimate: a well-designed prompt or embedding classifier might achieve **80–95% labeling accuracy** after tuning. The Amazon patent hints at top-model accuracy ~80% on return causes (Source: patents.google.com), although that was a more complex root-cause inference. Even 80% auto-labeling (with 20% needing review) would save enormous effort.

## Example Statistics

To frame the potential impact, consider an example retailer:

- **Case A: Mid-Sized Apparel Retailer.** Sells 100,000 items/month. With a 20% return rate (common for apparel), it processes ~20,000 returns. If reasons are free-text, manually tagging each is impractical. Automating at 80% accuracy means 16,000 auto-labeled, 4,000 flagged for review.

- **Manual Cost:** If an agent spends 2 minutes per return to read and categorize, labeling 20k takes ~667 person-hours (over 80 workdays). At $25/hour, that's $16,667 monthly!

- **AI Cost:** Using N/LLM embed (with free quota before hitting cost), 20k embeddings (1000 tokens each) is on the order of $ hundreds at current OCI pricing – far cheaper than manual labor. Even generative prompts (at, say, 20 tokens output + 100 tokens input per request) for 20k items might cost $ low-thousands monthly.

The ROI is clear: automating classification slashes labor *and* yields faster insights. Even if Amazon-like enterprises report an 80% predictive accuracy (Source: patents.google.com), this is acceptable if coupled with spot-checking.

## Comparative Studies

While no public study examines NetSuite specifically, we draw analogies. Academic work on complaint categorization (a similar NLP task) shows that multi-class classification with advanced models routinely surpasses rule- and keyword-based baselines (Source: ashwin-ks.github.io) (Source: houseblend.io). Moreover, embedding-based classification has been effective in real projects: e.g. LinkedIn data scientists have used pre-trained sentence embeddings to automatically tag user posts with thematic labels much faster than manual curation.

In product return literature, researchers have applied topic modeling on customer reviews to *infer* return causes (Source: link.springer.com). Their open model enriched the generic return code list – demonstrating that data-driven analysis reveals nuanced reasons. We propose to go further: not only analyze historical reasons, but also auto-apply standard labels going forward.

## Case Studies and Examples

Exact case studies on "NetSuite N/LLM for refund classification" are scarce (this is a cutting-edge use case). However, analogous examples illustrate the approach:

- **Internal Pilot (Hypothetical):** Consider a NetSuite user "TechGear Inc." using N/LLM. Initially, customer support logged return reasons in a free-text note. After implementing N/LLM classification, they defined 6 categories: *Defect, Fit/Size, Wrong Item, Late Delivery, Description Mismatch, Other*. Using a prompt-based approach, they tested 1,000 past returns and found the LLM assigned the expected category in ~85% of cases. In one quarter, automation reduced manual classification time by 75%, enabling weekly summary reports trending defect rates by product line. (Metric: previously, engineering saw defect spikes **1 month late**; now immediate alerts flow from AI labels.)

- **BirdRock Home (Real Example):** While not a refund classification study, BirdRock Home – a NetSuite customer – exemplifies a business fit for this technology. BirdRock has thousands of products and "processes thousands of orders daily" within NetSuite (Source: houseblend.io). Their unified data allowed them to build predictive churn models and inventory forecasts (Source: houseblend.io). One can imagine a similar workflow: BirdRock collects all return comments for its home-goods items and runs them through N/LLM. Because their catalog items have attributes (e.g. product descriptions) already in NetSuite, they could RAG-augment each reason with the item's details, further improving accuracy. For example, a complaint on a "sofa" about "sagging cushions" might be classified as *Item Not as Described*. The NetSuite platform thus becomes an AI-rich hub: BirdRock's example suggests large-scale returns processing is already part of their support volume, primed for automation.

- **Hypothetical Multi-national Retailer:** Imagine a retailer uses NetSuite globally. Return processes differ by country (e.g. different language, policy regimes). N/LLM's multilingual capabilities could handle varied languages (embedded translator module) and prompt templates. For instance, a Spanish-language reason "Demasiado pequeño" could be auto-detected as *Sizing*. This showcases the advantage of an ERP-native LLM: notifications and data stay within the secure enterprise environment, respecting localization settings (Source: houseblend.io) (Source: docs.oracle.com).

These scenarios underscore that N/LLM can integrate into real workflows: running as scheduled SuiteScripts, updating NetSuite records, and providing near-cahoots analysis. They also highlight the potential pitfalls: BirdRock's example shows "thousands of orders daily" – effective LLM use must scale to such volume and guard against overuse (monitoring quotas and cost).

## Implications and Future Directions

### Business Impact

Standardizing refund reasons with AI has multiple benefits:

- **Improved Analytics:** Management reports can break down returns by cause with confidence. For example, finance could see "25% of returns due to sizing/kitting issues this quarter, up from 20% last quarter," or logistics may note "defects concentrated in one product batch." Such insights support targeted improvements (e.g. adjusting size charts, improving quality control).

- **Fraud Detection and Policy Enforcement:** With semantic labels, unusual patterns (e.g. an item with a suspiciously high "wrong item" rate) can raise flags. Automating the initial review of cases against policy (e.g. excessive "didn't arrive" claims) becomes easier when reasons are normalized.

- **Customer Experience:** Agents can retrieve standardized reasons internally, ensuring consistent language in communications.

- **Cost Savings:** As estimated, reducing manual review by even a fraction yields labor savings. Aggregated, AI-driven returns processing may cut support staffing needs by tens of percent. For example, Topmost Labs reports that retailers deploying refund triage AI have halved their support costs (Source: topmostlabs.com) (though exact figures depend on workflow).

## Governance, Risk, and Compliance

Embedding AI within NetSuite raises governance issues, as highlighted by analysts. The Palma.ai review warns that **without governance, AI projects fail** (Source: palma.ai). In our context, we must ensure:

- **Data Privacy:** Refund reasons may include sensitive info ("address illegible", "child safety harness"). Any AI use should comply with privacy rules. NetSuite's design helps: N/LLM calls stay within the Oracle cloud and under NetSuite security controls. But developers should still ensure PII (names, etc.) are redacted or treated appropriately.

- **Audit Trails:** NetSuite logs changes to records. When a script writes the auto-assigned category, that change is recorded. Additionally, N/LLM can be set to return citations indicating which documents influenced a response (Source: docs.oracle.com), aiding transparency. All prompt inputs/outputs should be logged if required by policy.

- **Accuracy Monitoring:** As RecordPoint notes, **classification automation reduces human error but introduces ML risk** (Source: www.recordpoint.com). Organizations should routinely sample and verify labels, adjusting prompts or thresholds as needed. Keeping a small proportion of returns for human double-check (and feeding back corrections) creates a feedback loop. Over time, the AI model "learns" the company's vocabulary.

- **Regulatory Compliance:** In regulated industries (finance, pharmaceuticals), the reasons for returns (or refunds in payments) might have legal significance. One must ensure the AI categorization itself does not distort regulatory reporting. For example, if a recall occurs, the AI output should not mask this by labeling it simply as "Defective" without noting context. In certain cases, it may be prudent to keep the original free-text or concepts accessible for auditors.

Fortunately, NetSuite's approach – linking AI with in-system workflows – aligns with best practices: it avoids exposing data to ungoverned third-party AI tools (no "shadow AI" as Palma describes (Source: palma.ai). The Built-in SuiteScript AI modules operate under the account's security realm, requiring explicit enabling of the "Generative AI" feature. Still, businesses should define usage policies (who can run classification scripts, how often, etc.) and ensure that N/LLM usage logs are reviewed (e.g. limit number of API calls or data fields accessible).

## Future Directions

Looking ahead, several extensions are plausible:

- **Sophisticated Category Models:** Beyond flat labels, one could build multi-label or hierarchical classification (e.g. subcategories of "Defect": cosmetic vs functional). N/LLM could even generate short explanations ("Broken screen assembly").

- **Language Support:** As noted, NetSuite's language modules can translate return reasons. Multi-lingual classification could be achieved by first translating reasons to a pivot language or by using GPT-4-level models.

- **Continuous Learning:** Rather than fixed prompts, the system could automatically refine its categories. For example, the recordpoint article suggests AI classifiers "continuously improve their accuracy by learning from classification decisions and feedback loops" (Source: www.recordpoint.com). In NetSuite, this might look like periodically retraining on any confirmed corrections.

- **Integration with Workflow:** Classified reasons could drive automated actions. For instance, any "Defective" return might automatically trigger a vendor quality check or refund priority; any "Wrong Item" could auto-issue free return shipping.

HOUSEBLEND

- **LLM Assistant Interfaces:** With NetSuite's planned "SuiteAnalytics Assistant" (natural-language query of data (Source: houseblend.io), users might eventually get charts by just asking "Show refunds by reason for last quarter." The accuracy of such assistants will depend on proper classification upstream.

- **Benchmarking and Best Practices:** As N/LLM (and generative features like Text Enhance (Source: www.techtarget.com) proliferate, NetSuite will likely publish more guides or even built-in flows for use cases. By late 2025, we anticipate Oracle or partners providing "Refund Reason Classification" as a documented reference scenario, complete with recommended prompts.

## Conclusion

Standardizing refund reasons is a high-value problem for any business dealing with product returns or payment reversals. Inconsistent or free-form reason data obscures root causes and inflates costs. The recent advent of enterprise LLM capabilities – notably Oracle NetSuite's N/LLM module – offers a powerful solution. Using embeddings and generative AI directly within the NetSuite environment, companies can automatically categorize refund comments at scale.

This report has examined the context (scope of returns, common causes), the challenges of unstructured reason data, and how N/LLM technical features address these challenges. We have drawn on industry data (e.g. **42% of returns due to fit issues** (Source: corp.narvar.com), **$890B in U.S. returns** (Source: www.glencoyne.com) to underscore the opportunity. We detailed a classification workflow: extracting reason text via SuiteQL, defining category taxonomy, invoking LLM embeddings or prompts, and writing back standardized labels.

As NetSuite and Oracle strengthen their AI roadmap (e.g. suitewide Text Enhance features (Source: www.techtarget.com), we expect such text-classification applications to become more turnkey. Crucially, success depends on sound data governance and validation (Source: palma.ai) (Source: www.recordpoint.com). With proper oversight, however, the ROI is compelling: faster analytics, better decision-making, and major savings in manual effort.

In sum, **auto-classification of refund reasons via NetSuite N/LLM** exemplifies the next wave of ERP+AI: leveraging native AI modules to turn unstructured data into actionable intelligence. The early signals — from technical docs and initial pilots — suggest organizations that embrace this will gain sharper insights into their returns processes and unlock hidden value in their transactional data.

## References

- Oracle NetSuite SuiteScript N/LLM Module documentation (Source: docs.oracle.com) (Source: docs.oracle.com) (Source: docs.oracle.com)
- Houseblend, *"Enriching NetSuite Data: A Guide to the N/LLM Module" (Nov 2025)* (Source: houseblend.io) (Source: houseblend.io)
- TechTarget, *"Oracle NetSuite follows generative AI trend in ERP"* (Oct 2023) (Source: www.techtarget.com) (Source: www.techtarget.com)
- Glencoyne, *"True Cost Per Return: Hidden Cash Drain"* (Jun 2025) (Source: www.glencoyne.com) (Source: www.glencoyne.com)
- Narvar (by Narvar Team), *"6 Common Retail Return Reasons"* (Jul 2025) (Source: corp.narvar.com) (Source: corp.narvar.com)
- Springer (E-Commerce Research), *Zangara et al., "Shaping the causes of product returns: topic modeling…"* (Sep 2024) (Source: link.springer.com) (Source: link.springer.com)
- US11526665B1 (Amazon Technologies) – Patent on ML root-cause of returns (Source: patents.google.com)
- Reuters / DataIntelo, *"Refund Abuse Detection AI Market Report"* (2023) – Industry statistics on e-commerce returns. (Note: for context)
- Palma.ai (blog), *"Why Enterprise AI Fails Without Proper Data Governance"* (Source: palma.ai)
- RecordPoint, *"How AI classification solves information governance challenges"* (Source: www.recordpoint.com)
- Oracle NetSuite Help Center – Return Authorization and Customer Returns documentation (Source: docs.oracle.com) (Source: docs.oracle.com)
- Oracle NetSuite Community (user threads on return processing) – for procedural insights.
- Industry analyses (e.g. Narvar, Techradar, Axios) on ERP-AI trends (Source: houseblend.io) (Source: www.techtarget.com).

Tags: netsuite n/llm, generative ai, suitescript, text classification, refund analysis, data standardization, large language models, e-commerce returns