# Making Enterprise AI Stick:

## How Solving Two Adoption Failures Transformed Salesforce Service Cloud AI

Anjan Goswami

### Abstract

Salesforce Service Cloud AI shipped two AI-powered capabilities for customer service: a chatbot that automated routine customer interactions through intent classification and templated responses, and an answer recommendation engine that surfaced relevant knowledge base articles and forum answers to human agents in real time. Both were technically functional. Neither was being adopted. This case study describes how diagnosing two distinct adoption failures— customers could not configure their chatbots, and agents did not trust the answer suggestions— led to fundamentally different interventions: vertical-specific defaults with purchased utterance data for the chatbot, and a dual-index retrieval architecture with aggressive noise suppression for the answer engine. The chatbot fix was a product design insight (invert the setup burden so customers provide answers, not utterances). The answer engine fix was an engineering architecture insight (separate curated content from noisy forum content at the index level, and treat precision as more important than recall). Together they transformed Service Cloud AI adoption and the system remained the production architecture for approximately five years, until LLM-based approaches superseded the intent classification paradigm.

## 1 The Mandate: Make Service Cloud AI Successful

In 2018, Salesforce Service Cloud AI had two flagship capabilities. The **chatbot** automated customer service conversations: customers interacted with a bot that classified their question into an intent, collected relevant information through a multi-turn dialogue, and returned a templated answer. The **answer recommendation engine** helped human agents during live customer interactions by surfacing relevant articles and answer snippets in real time as the customer described their problem.

Both systems worked in controlled demonstrations. Neither was succeeding in the field. I was brought in to diagnose why and fix it.

The diagnosis revealed that the two systems were failing for **completely different reasons**— and therefore required completely different interventions. Understanding this distinction was the critical first step. A single "improve the AI" initiative would have addressed neither problem. What was needed was a precise diagnosis of each adoption failure, followed by a targeted fix matched to the root cause.

## 2 Adoption Failure #1: Customers Could Not Launch Their Chatbots

### 2.1 How the Chatbot Worked

The chatbot used a pre-LLM conversational AI architecture: a state-preserving dialogue system driven by intent classification. The customer was greeted with a templated response, provided identifying information (name, account number), asked their question, and the system classified the question into one of a predefined set of intents—"check order status," "reset password," "billing dispute." Each intent had a templated answer with dynamic fields: "Your order #{order_id}

shipped on {ship_date} via {carrier}." Multi-turn conversations were handled through a state machine tracking dialogue position and collected information.

The intent classifier was BERT-based, later extended to **multilingual BERT** supporting 11–20 languages. The multilingual architecture was essential for Salesforce's global customer base—a Japanese retailer, a German bank, and a Brazilian telecom company all needed chatbots in their local languages. Multilingual BERT's cross-lingual transfer allowed training primarily on English utterance data with minimal per-language fine-tuning, handling code-switching and multilingual queries from a single model.

The technology worked. **The setup process did not.**

## 2.2   The Setup Problem: Customers Did Not Know How to Write Utterances

Salesforce shipped a training interface where customers populated their chatbot by providing intents, example utterances for each intent, and templated answers. Step 2—writing utterances—was the adoption killer.

Writing good training utterances requires understanding what a classifier needs: diverse phrasings, edge cases, near-miss examples that distinguish similar intents. A customer service manager at a retail company knows that customers ask about order status. They do not know that the classifier needs 50–100 diverse phrasings including "track my package," "when will my stuff arrive," "order hasn't come yet," "shipping update please," and "did my order ship."

Customers provided 5–10 utterances per intent with minimal lexical diversity. The classifiers trained on this data were unreliable. Chatbots misclassified intents, returned wrong answers, and were abandoned.

Even motivated customers with adequate utterances faced a deeper problem: **per-organization training produced weak classifiers.** A single organization might have 20–50 intents with 10–20 utterances each—a few hundred training examples total. BERT fine-tuning on this volume overfit rapidly and failed to generalize to the natural language variation in real customer conversations.

## 2.3   What I Tried First That Failed

My initial approach was to improve the training interface—better documentation, guided wizards for utterance creation, example utterances that customers could copy and modify. **This did not move the needle.** The problem was not UX friction in the interface. The problem was that customers fundamentally could not produce the data the system needed. A better form for entering utterances did not help people who did not know what utterances to enter.

I also explored automatic utterance augmentation—paraphrasing customer-provided utterances using rule-based transformations (synonym substitution, word reordering, template expansion). The augmented data improved classifier robustness marginally but could not compensate for the narrow semantic coverage of the original utterances. If a customer only provided "where is my order" and close variants, augmentation produced more variants of that phrasing but never discovered "track my package" or "shipping update please."

Both attempts failed because they tried to make the existing process work better. The process itself was wrong.

## 2.4   The Fix: Vertical-Specific Defaults

The breakthrough came from recognizing that **customer service questions follow strong vertical patterns.** "Where is my order?" is asked at every retailer. "What's my account balance?" is asked at every bank. "How do I file a claim?" is asked at every insurance company. The utterances vary in phrasing but the intents are universal within verticals.

I purchased common utterance datasets from specialized vendors—curated collections of customer service questions organized by industry vertical, with hundreds of phrasings per intent covering colloquial language, misspellings, abbreviations, and multilingual variations.

Using this data supplemented by Salesforce's own customer service data, I constructed **vertical-specific default configurations**:

| Vertical | Example Default Intents | Intent Count |
|---|---|---|
| Retail / E-commerce | Order status, returns, shipping, product info | 30–50 |
| Banking / Finance | Balance inquiry, transactions, card issues | 25–40 |
| Insurance | Claims filing, policy questions, renewals | 20–35 |
| Telecommunications | Service outage, billing, plan changes | 25–45 |
| Technology / SaaS | Password reset, feature help, bug reports | 20–40 |

Table 1: Vertical-specific default intent configurations.

Each vertical template shipped with pre-trained intent classifiers (already trained on hundreds of utterances per intent) and **blank answer templates**. The customer's onboarding was inverted:

| Before (Failed) | After (Vertical Defaults) |
|---|---|
| 1. Customer defines intents from scratch | 1. Customer selects industry vertical |
| 2. Customer writes 50–100 utterances per intent | 2. System loads pre-trained intents |
| 3. Customer writes answer templates | 3. Customer fills in answers (fill-in-the-blank) |
| 4. Classifier trains on sparse data | 4. Classifier already trained; optional fine-tuning |
| 5. Chatbot launches with poor accuracy | 5. Chatbot launches with production-quality accuracy |

Table 2: Onboarding before and after vertical defaults.

The core insight: **customers know their answers but not their utterances.** A customer service manager can immediately write "Your refund will be processed within 5–7 business days and credited to your original payment method." They cannot write 100 ways a customer might ask for a refund. By providing the hard part (utterance diversity) and asking customers for the easy part (answers), we removed the primary adoption barrier.

Per-organization fine-tuning used the vertical default model as initialization, with any customer-provided additional utterances applied as a brief fine-tuning pass. Transfer learning from a well-trained default produced substantially better classifiers than training from scratch, even with minimal per-customer data.

# 3 Adoption Failure #2: Agents Did Not Trust the Answer Suggestions

## 3.1 How the Answer Engine Worked

Many customer interactions required human agents—complex issues, escalations, edge cases beyond the chatbot's intent coverage. The answer recommendation engine surfaced relevant articles and answer snippets to agents in real time, drawn from two sources: curated knowledge base articles and historical forum Q&A pairs.

When the chatbot could not classify a customer's intent with sufficient confidence, the conversation escalated to a human agent. The answer engine received the conversation history and immediately began surfacing suggestions—providing continuity between the automated and human-assisted phases.

## 3.2  Why Agents Stopped Using It

The first version of the answer engine used standard BM25 retrieval against a single combined index of articles and forum content. Agents initially tried the suggestions, then stopped. I investigated why.

The problem was **noise from forum content contaminating the results.** Forum Q&A threads accumulated dozens of responses per question, and most responses were not answers:

- Comments that critiqued the question rather than answering it

- Follow-up questions from other users with tangentially related problems

- Outdated solutions referencing deprecated product versions

- Tangential discussions between respondents

- Answers that were simply wrong

BM25 matched these responses on lexical overlap with the customer's question and ranked them alongside curated articles. An agent who received a wrong forum answer and sent it to a customer learned not to trust the system. **A few bad suggestions destroyed trust faster than many good suggestions built it.**

The second problem was **stale content.** BM25 is recency-agnostic. An article written three years ago about a deprecated product version ranked highly if it matched the query terms. Agents who followed outdated instructions compounded the trust problem.

## 3.3  What I Tried First That Failed

My first fix was post-retrieval filtering: retrieve from the combined index, then apply quality heuristics to filter out low-quality results before showing them to agents. This improved results slightly but had a fundamental flaw—**the retrieval stage itself was polluted.** When noisy forum content occupied retrieval slots, high-quality articles that should have been retrieved were pushed below the retrieval cutoff. Filtering after retrieval could remove bad results but could not recover good results that were never retrieved.

I also tried boosting curated articles in the combined index by multiplying their BM25 scores by a source-quality factor. This helped but created a new problem: the boost factor had to be set globally, which over-weighted mediocre articles relative to genuinely excellent forum answers. Some forum Q&A pairs were the best answers available—written by expert community members with deep product knowledge. A global source boost suppressed these alongside the noise.

## 3.4  The Fix: Dual-Index Architecture with Learned Ranking

The solution required structural separation, not parameter tuning.

**Separate indexes**: I split the single index into two: one for curated knowledge base articles and one for forum Q&A content. For the forum index, Q&A threads were **flattened into question-answer pairs**—a thread with one question and five responses became five separate retrieval documents. This enabled direct matching between a customer's current question and historical question-answer pairs, bypassing the noisy thread structure.

The separation allowed **source-specific retrieval strategies**: different quality priors, different scoring adjustments, different filtering thresholds for each source. Curated articles had an inherent quality prior from editorial review. Forum content required aggressive noise suppression.

**Beyond BM25**: I augmented retrieval with data quality factors:

$$\text{Score}(q, d) = \text{BM25}(q, d) \cdot \phi_{\text{quality}}(d) \cdot \phi_{\text{recency}}(d) \cdot \phi_{\text{source}}(d)$$

The quality factor $\phi_{\text{quality}}(d)$ incorporated article curation status, forum answer engagement signals (upvotes, acceptance marks), author credibility, and critically, a **noise classifier** trained

to distinguish actual answers from critiques, clarifying questions, and off-topic discussion. The recency factor $\phi_{\text{recency}}(d)$ applied different decay rates: slower for curated articles (which age more gracefully) and faster for forum answers (which become outdated quickly). The source factor $\phi_{\text{source}}(d)$ provided a baseline quality advantage to curated articles—not a global override, but a prior that the learned ranker could adjust per query.

**The learned ranking model**: On top of retrieval, I trained a learning-to-rank model using human relevance judges. I recruited evaluators who rated answer candidates on a graded relevance scale (0–4) given the customer's question. The initial model used approximately 20 features:

- BM25 scores from both indexes

- Quality factors (curation status, engagement signals, author credibility)

- Textual similarity features (TF-IDF cosine, Jaccard overlap, query term coverage)

- Document metadata (article type, creation date, last-updated, view count)

- Source type and noise classifier confidence

The relevance judges served double duty: validating overall system quality and providing training signal for the ranker. As the system accumulated behavioral feedback from agent interactions—which suggestions agents selected, modified, or ignored—the feature set expanded from **20 features to over a thousand**, incorporating fine-grained textual features, cross-field matching signals, agent behavioral patterns, and organization-specific relevance signals.

## 3.5 Optimized for Small Collections

A critical design choice was recognizing that **most organizations had fewer than 200,000 documents**. This constraint was an asset:

- The entire index fit in memory, enabling sub-millisecond retrieval during live agent conversations

- The small candidate set allowed applying the full ranking model to all retrieved documents without cascading, eliminating information loss

- Hardware requirements were derived directly from collection size, enabling precise per-customer capacity planning

## 3.6 How Agents Responded

After the dual-index rebuild, agent engagement with the answer engine recovered. The key metric was not click-through rate on suggestions—it was **whether agents continued using the system over time.** With the original system, usage decayed week over week as agents lost trust. With the rebuilt system, usage was sustained.

Agents reported that the most valued behavior was **what the system did not show them.** The noise suppression—filtering critiques, outdated content, and wrong answers—was invisible when working correctly, but its absence had been the original system's downfall.

# 4 The Two Fixes in Retrospect

The two problems required fundamentally different thinking. The chatbot problem was a **product design failure**—the system asked customers to do something they could not do, and the fix was inverting the burden. The answer engine problem was an **engineering architecture failure**—the system mixed clean and noisy data at the structural level, and the fix was separating them at the index level rather than filtering after the damage was done.

|  | **Chatbot** | **Answer Engine** |
|---|---|---|
| **Symptom** | Low adoption; customers abandon setup | Low usage; agents stop trusting suggestions |
| **Root cause** | Customers cannot provide training data the model needs | Noisy forum content and stale articles erode trust |
| **Failed fix** | Better training UI; utterance augmentation | Post-retrieval filtering; global source boosting |
| **Why it failed** | Wrong assumption: customers can produce utterances if we help | Wrong architecture: filtering after retrieval cannot recover displaced results |
| **Successful fix** | Vertical defaults: provide utterances, ask for answers | Dual-index + noise suppression + learned ranking |
| **Nature of fix** | Product design insight | Engineering architecture insight |

Table 3: Comparison of the two adoption failures and their resolutions.

Recognizing which type of problem you are facing—product design vs. engineering architecture—determines whether you iterate on UX and data strategy or rebuild the technical foundation. Attempting the wrong type of fix wastes months.

# 5 Results

| **Metric** | **Result** |
|---|---|
| Chatbot adoption | Significant increase after vertical defaults launch |
| Chatbot intent accuracy | Production-quality from day one (vertical defaults) |
| Language coverage | 11–20 languages via multilingual BERT |
| Answer engine agent usage | Sustained engagement (vs. decaying usage pre-fix) |
| Answer engine handle time | Reduced average agent handle time |
| Ranking model scale | Grew from 20 to 1,000+ features over system lifetime |
| System longevity | Production architecture ∼5 years (2018–2023) |

Table 4: Service Cloud AI impact after both adoption fixes.

The system remained the production architecture for approximately five years. The answer engine's retrieval-and-ranking infrastructure proved particularly durable: even as LLMs replaced the chatbot's intent classification stack beginning in 2023, the retrieval indexes, quality scoring, and ranking models continued to serve as the knowledge grounding layer for LLM-generated responses. The infrastructure built to solve a pre-LLM problem became the RAG pipeline for the post-LLM era.

# 6 Generalizable Principles

**1. Diagnose which type of adoption failure you have before fixing anything.** The chatbot and answer engine failed for completely different reasons. A single "improve the AI" initiative

would have addressed neither. Product design failures require inverting user burden. Engineering architecture failures require structural rebuilds. Applying the wrong fix type wastes months.

**2. Customers know their answers, not their questions.** This asymmetry is fundamental in enterprise AI. Domain experts can articulate solutions but cannot enumerate the ways users phrase problems. Any system requiring domain experts to provide question-side training data will face adoption resistance. Invert the burden: provide the question variations, ask for the answers.

**3. A few bad suggestions destroy trust faster than many good suggestions build it.** In agent-assist systems, agents will tolerate "no suggestions available" but will permanently abandon a system that serves them wrong answers. Precision dominates recall. Invest in noise suppression and data quality before optimizing for coverage.

**4. Separate noisy sources from curated sources at the index level.** Post-retrieval filtering cannot recover good documents displaced during retrieval by noisy ones. Structural separation at the index level—with source-specific quality priors, decay rates, and scoring—outperforms any amount of post-hoc filtering on a combined index.

**5. Vertical defaults solve the per-organization cold-start problem.** Individual organizations have too little data to train robust models. But across organizations within a vertical, the same patterns recur. Purchasing or aggregating cross-organization data and packaging it as vertical defaults transforms a cold-start problem into a transfer learning problem.

**6. Retrieval infrastructure outlasts the generation paradigm.** Intent classification was replaced by language models. The retrieval indexes, quality scoring, and ranking models survived the transition and became the RAG pipeline. Building clean retrieval and ranking infrastructure is a durable investment regardless of what sits on top of it.

**7. The product is not the model—it is the model plus the setup experience.** A technically superior system that customers cannot configure produces zero value. The vertical defaults did not improve the underlying BERT classifier. They removed the barrier between the classifier and the customer. In enterprise AI, the setup experience is as much a part of the product as the model itself.

---

Anjan Goswami, Ph.D.                                    anjangoswami.com    ·    smartinfer.com