



---

# AI営業電話品質評価レポート2026 - STT・LLM・TTS統合品質の技術分析 | Leadsia Inc.

2026-03-25

Leadsia Inc.

## AI営業電話の品質を決める3つの要素 | なぜ人間レベルの対話が可能なのか

### カテゴリー：AI導入の基礎知識

「AI営業電話」と聞いて、多くの方がイメージするのはこんな光景ではないでしょうか。

機械的な声が一方向的に話し続け、相手が何か言おうとしても無視して次のセリフに進む。途中で質問しても「もう一度お願いします」と繰り返す。不自然な間が空いて、相手に「これ、ロボットだな」とバレル。

確かに、そういうAI電話は今でも存在します。

しかし2026年現在、音声AI技術は「人間と区別がつかない」レベルの対話品質を実現し始めています。ベテラン営業マンのような滑らかな受け答え、相手の割り込みへの自然な対応、文脈を踏まえた的確な質問 - これらは、もはやSFの話ではありません。

では、「不自然なAI電話」と「人間レベルのAI電話」を分けるものは何か。本記事では、AI営業電話の品質を決める3つの要素を解説します。

#### 要素①：「聞く力」 - 日本語音声認識の精度が会話の土台

AI営業電話の品質は、まず「相手の声をどれだけ正確に聞き取れるか」で決まります。

これを担うのが音声認識（STT=Speech-to-Text）技術です。電話口の音声をリアルタイムでテキストに変換し、AIの「頭脳」に渡す。この変換の精度が低ければ、どれだけ優秀なAIを搭載しても、会話は成り立ちません。

日本語は特に難しい

音声認識において、日本語は世界でも難度の高い言語の一つです。

まず、同音異義語が非常に多い。「きかん」だけでも「期間」「機関」「器官」「帰還」「既刊」 - 文脈なしでは正解を特定できません。ビジネス電話では「貴社」と「記者」、「ご請求」と「ご精究」のような聞き間違いが致命的です。

次に、敬語表現の複雑さ。「お伺いしたいのですが」「ご確認いただけますでしょうか」といったビジネス日

本語は、層の深い敬語構造を持ちます。これを正確に認識できなければ、AIの応答も不自然になります。

さらに、電話特有のノイズ環境。オフィスの雑音、携帯電話の通話品質、相手の話し方の個人差 - これらすべてが認識精度を下げる要因です。

#### 汎用モデルと日本語特化モデルの差

海外で開発された汎用の音声認識モデルは、英語では高い精度を発揮しますが、日本語ではその性能を十分に発揮できないケースがあります。

例えば、医療分野での比較検証では、日本語に特化した音声認識エンジンが認識精度97%超を記録したのに対し、汎用モデルは84~87%程度にとどまったという報告があります。約10ポイントの差は、100語の発話で10語以上の聞き間違いが生じることを意味します。営業電話でこれは致命的です。

国産の日本語特化型エンジンは、ビジネス用語、業界専門用語、方言（関西弁や東北弁など）への対応力で優位性を持っています。また、コールセンターのような騒がしい環境でのノイズ除去技術も、日本国内の通話環境に最適化されています。

AI営業電話の品質を左右する最初の分岐点は、「日本語の聞き取り精度」です。

#### 要素②：「考える力」 - 大規模言語モデル（LLM）が会話の質を決める

音声認識で正確にテキスト化された相手の発話を、「理解」し、「適切な応答を生成」する。これを担うのが大規模言語モデル（LLM） - AI営業電話の頭脳部分です。

#### キーワードマッチングとLLMの決定的な違い

従来のAI電話システムの多くは、「キーワードマッチング」で応答を選択していました。相手の発話から「見積もり」「料金」「日程」などのキーワードを検出し、あらかじめ用意された定型回答を返す仕組みです。

この方式は、「見積もりをください」のような単純な発話には対応できます。しかし、こんな場面ではどうでしょうか。

「先週もらった見積もり、ちょっと上の者と相談したんですが、もう少しだけ条件面で融通きかないかなって……あ、でも予算的にはほぼ固まってるんで、大幅な変更は考えてないです」

この発話には、「再交渉の打診」「上長の関与」「予算は確保済み」「大幅値引きは不要」という複数の情報が含まれています。キーワードマッチングでは「見積もり」と「条件」を拾って「新しい見積もりを作成します」と返すのが精一杯。しかし、ベテラン営業マンなら「条件面というのは、具体的にどのあたりをご検討でしょうか？」と自然に掘り下げますよね。

LLMは、この「文脈全体を理解した上での適切な応答生成」を実現します。単語の羅列ではなく、発話の意図・背景・感情を総合的に判断し、会話の流れに沿った自然な返答を生成する。

すべてのLLMが同じではない

ただし、LLMの品質は一様ではありません。営業電話という用途では、特に重要な特性があります。

誠実さ。相手に媚びて実現不可能な約束をしないか。「なんでもお安くします！」と勝手に値引き交渉をしないか。事実と異なることを自信満々に言わないか。

一貫性。会話の途中で矛盾する発言をしないか。前半で言ったことと後半で言ったことが食い違わないか。

安全性。不適切な発言をしないか。相手の感情を不用意に増幅させないか。競合製品について事実に基づかない批判をしないか。

2025年にOpenAIのGPT-4oで問題となった「sycophancy（追従性）」 - AIがユーザーの意見に過剰に同調する傾向 - は、このLLMの「性格」に起因する問題でした。

LeadsiaのAI営業インテリジェンス「ALICE」が採用しているAnthropicのClaudeは、「Constitutional AI（憲法AI）」というアプローチで訓練されています。明文化された原則に基づいてAI自身が応答を評価・修正する仕組みにより、おべっかを言わず、事実に基づいた誠実な会話を設計レベルで担保しています。

（詳しくは「AIのモデルに性格はあるのか？」をご参照ください）

AI営業電話の品質の核心は、「どのLLMを、どう活用しているか」にあります。

要素③：「話す力」 - 音声合成（TTS）が第一印象を決める

相手の声を聞き取り、適切な応答を考えた。最後に、その応答を「自然な音声」で伝える。これを担うのが音声合成（TTS=Text-to-Speech）技術です。

TTSは「読み上げ」から「演技」のフェーズへ

数年前のTTSは、明らかに「ロボットの声」でした。抑揚が不自然で、どこか冷たい印象。電話口で聞いた瞬間に「あ、これ自動音声だ」とわかるレベルでした。

2026年現在、TTS技術は劇的に進化しています。

最新の音声合成エンジンは、吐息、笑い、微細な感情の揺れまで再現可能になっています。話速の自然な変化、間（ま）の取り方、共感を示すトーン - 人間の声優と区別がつかないレベルに到達しつつあります。

業界最高水準のTTSエンジンでは、発話品質の誤認率（WER）が2.83%にまで低下しており、ほぼ完璧な自然言語を音声で再現できるようになっています。またストリーミング遅延も200ミリ秒程度にまで短縮され、リアルタイム対話に十分な応答速度を実現しています。

営業電話に求められるTTSの条件

しかし、「聞き取りやすい声」であるだけでは営業電話には不十分です。

営業の現場では、相手の反応に合わせてトーンを変える場面が頻繁にあります。興味を示してくれたら少し明るい声で。慎重な態度なら落ち着いたトーンで。断られそうなら共感を示す声色で。

最新の調査によると、感情検知機能を備えた音声AIは、顧客の苛立ちを検知して適切にエスカレーション（人間への引き継ぎ）を行うことで、クレーム発生率を約25%低減できるとされています。音声合成における「感情知能（Emotional Intelligence）」は、顧客体験を左右する重要な要素になっているのです。

AI営業電話の第一印象は、TTSの品質で決まります。

3つの要素が「統合」されて初めて品質が生まれる

ここまで、音声認識（聞く力）・LLM（考える力）・音声合成（話す力）の3要素を個別に解説してきました。

しかし、最も重要なポイントはこの3つが高いレベルで「統合」されているかどうかです。

音声認識が優秀でも、LLMの応答品質が低ければ的外れな返答になります。LLMが優秀でも、TTSが不自然なら「ロボットと話している」という印象を与えます。TTSが美しくても、音声認識の精度が低ければ会話が噛み合いません。

さらに、3つの要素を統合する際に生じる「レイテンシー（遅延）」の問題があります。

相手が話し終わってからAIが応答するまでの時間。これが長いと、電話口で不自然な沈黙が生まれます。人間の会話では、相手の発話終了から応答開始まで通常200~500ミリ秒。これより大幅に遅いと、相手は「聞こえていますか？」と不安になります。

音声認識→テキスト変換→LLM処理→応答生成→音声合成→音声出力 - この一連のパイプラインを、人間の会話テンポに合わせてリアルタイムで処理する。これは技術的に極めて高度な課題であり、この統合品質こそが「音声AIインテリジェンス」と従来のAI電話の決定的な差です。

（「音声AIインテリジェンス」の技術カテゴリーについては「音声AIインテリジェンス技術とは？」で詳しく解説しています）

「割り込み対応」 - 自然な会話を決定づける技術

3要素の統合品質を最もわかりやすく示すのが、「バージイン（割り込み対応）」です。

バージインとは、AIが話している最中に相手が発話した際、AIが即座に話をやめて「聞き取りモード」に切り替わる技術です。

人間の会話では当たり前のことですが、AI電話ではこれが非常に難しい。AIが自分の声を出しているとき、マイクにはAI自身の音声と相手の声が同時に入力されます。この中から「相手が話し始めた」ことを検知するには、高度な音声区間検出（VAD）とエコーキャンセラー技術が必要です。

バージインができないAI電話は、相手がどれだけ割り込もうとしてもAIが話し続けます。あるいは、少しの雑音でAIが黙り込んでしまい、会話が途切れます。どちらも、電話の相手にとってはストレスフルな体験です。

日本語音声AIにおけるバージインの実装事例は増えつつあります。自動車の音声ナビゲーションでは、案内音声の途中で「次の角を右！」と割り込める設計が安全性の観点から必須技術となっています。一部の国内銀行や保険会社のAI窓口でも、顧客が「あ、ちょっと待って」と言った瞬間にAIが沈黙し、次の発話を待つ自然なターンテイクングを実現しています。

ALICEは、この割り込み対応を営業電話の文脈で実現しています。相手が話し始めたら即座に聞く側に回り、相手の発話が終わったら適切なタイミングで応答を再開する。人間同士の会話と同じ「キャッチボール」が、AI営業電話で可能になっています。

## 品質を見極めるための実践的チェックリスト

AI営業電話の導入を検討する際、品質を見極めるために確認すべきポイントをまとめます。

デモ体験時にチェックすべきこと：

「聞く力」の確認として、早口で話しても正確に聞き取れるか、業界用語や固有名詞を認識できるか、周囲に多少の雑音がある状態でも機能するかを確認してください。

「考える力」の確認として、想定外の質問にどう応答するか、話題が途中で変わっても文脈を追えるか、曖昧な表現（「ちょっと考えさせてください」等）の意図を正しく理解するかを見てください。

「話す力」の確認として、声の自然さは十分か、間（ま）の取り方に違和感はないか、相槌のタイミングは自然かを確認してください。

そして「統合品質」として、相手が話し始めたらAIがすぐに黙るか（バージン対応）、応答までの沈黙時間は許容範囲か、会話全体を通じて「人間と話している」感覚があるかを総合的に評価してください。

これらすべてを高いレベルで実現しているサービスが、真の「音声AIインテリジェンス」です。

まとめ：「AIが話せる」と「AIと会話できる」は違う

AI営業電話の品質は、3つの要素で決まります。

聞く力（音声認識） - 日本語の精度が会話の土台。汎用モデルと日本語特化モデルでは10ポイント以上の精度差が生じることもあります。

考える力（LLM） - 文脈理解と応答生成がAIの「営業力」。モデルの「性格」（誠実さ、安全性）が営業品質を直接左右します。

話す力（音声合成） - 第一印象を決めるTTS品質。2026年現在、人間と区別がつかないレベルの技術が実用化されています。

そして、この3つが高いレベルで統合され、人間の会話テンポでリアルタイム処理される。割り込み対応、相槌、適切な間 - これらが揃って初めて、「AIと会話できる」体験が生まれます。

「AI営業電話」という言葉だけでは、その中身はわかりません。一方的に話すだけのシステムも、人間レベルの対話ができるシステムも、同じ名前と呼ばれています。導入前にデモを体験し、3つの要素と統合品質を自

分の耳で確かめること。それが、AI営業電話選びの最も確実な方法です。

#### 関連記事

- 音声AIインテリジェンス技術とは？従来の音声システムとの決定的な違い
- AIのモデルに「性格」はあるのか？ - Claudeの魂をつくった哲学者と、AIの人格設計という新領域
- セールステックSaaSの選び方 | 失敗しない5つのポイント
- ゼロタッチ運用とは？人間の作業時間を限りなくゼロにする設計思想

Leadsiaは、AI営業インテリジェンス「ALICE」、AI音声インテリジェンス「SOPHIA」、AI業務インテリジェンス「LYDIA」を通じて、日本のB2B企業の営業DXを支援するセールステックSaaS企業です。各AIエージェントの頭脳にはAnthropicのClaudeを採用し、Constitutional AI（憲法AI）に裏打ちされた安全性と会話品質を両立した営業自動化を実現しています。

詳しくは[Leadsia公式サイト]をご覧ください。