

Pragmatics of Synthetic Deception: Speech Act Analysis and Psycholinguistic Manipulation in Deepfake Voice Scam Messages in Indonesia

Feska Ajepri

STAI Ma'arif Kalirejo Lampung, Indonesia

*Email: feska_ajepri@staimaarifkalirejo.ac.id

Submitted: 08/02/2026	<i>Abstract</i>
Accepted: 01/03/2026	<p>Voice-cloned deepfake scams have become a fast-growing vector of financial fraud in Indonesia, with reported losses reaching the trillions of rupiah and a documented surge exceeding 1,500% in incident volume between 2022 and 2023. Forensic and cybersecurity literature has mapped the technical signatures of synthetic audio, but comparatively little work has examined how these synthetic voices linguistically perform deception: what speech acts they enact and how their prosodic design exploits listener cognition. This study addresses that gap by analyzing eight Indonesian-language deepfake scam voice samples, reconstructed from publicly reported cases and verified incident patterns, through Searle's (1976) taxonomy of illocutionary acts combined with discourse-psycholinguistic analysis of pause structure and pitch (Fo) contour using Praat-based acoustic modeling. Results show that directive acts dominate the corpus (36% of 95 coded units), typically preceded by assertive acts that construct false authority, and that scam discourse follows a recurring five-phase structure: authority framing, crisis-trigger diction, a manipulative pre-directive pause, directive command, and closing reassurance, a sequence absent from ordinary conversational speech. The pre-directive pause (M = 980 ms) runs roughly three times longer than the conversational baseline (M = 340 ms) and functions as a cognitively engineered interval that suppresses critical evaluation before the financial demand is issued. These findings extend speech act theory into synthetic media contexts and show that deepfake persuasion is a deliberately engineered pragmatic structure, not just a visual or technical artifact. The study contributes a replicable analytic framework for forensic linguists, digital literacy educators, and policymakers seeking to counter AI-generated voice fraud.</p>
Published: 31/03/2026	<p>Keywords: forensic linguistics; speech act theory; deepfake voice scam; psycholinguistic manipulation; prosodic analysis</p>

© The Author(s). 2026 Open Access. This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0>

INTRODUCTION

Synthetic voice cloning has moved from a laboratory curiosity to an operational instrument of financial crime in Indonesia within a remarkably short span. Industry monitoring by PT Indonesia Digital Identity (VIDA) recorded a 1,550% increase in deepfake-related fraud cases between 2022 and 2023, and by January 2026 the Indonesia Anti-Scam Centre and the Financial Services Authority (OJK) had logged more than 432,000 digital fraud reports with losses reaching the trillions of rupiah, a substantial share attributed to AI-generated voice and

video impersonation. The most widely reported domestic case involved a Lampung-based perpetrator, identified by police only by the initials AMA, who cloned the voices and likenesses of President Prabowo Subianto, Vice President Gibran Rakabuming Raka, and Finance Minister Sri Mulyani to promise fictitious cash assistance of Rp50 million, collecting roughly Rp30 million from eleven victims before his arrest in January 2025. A second perpetrator using the same script template, identified as JS, was arrested the following month in Pringsewu after extracting Rp65 million from approximately one hundred victims. A separate scheme in East Java used a cloned voice of Governor Khofifah Indar Parawansa to advertise motorcycles at an implausible Rp500,000, netting Rp87 million from roughly one hundred victims before three suspects were apprehended by the East Java Regional Police.

These cases share a structural feature that has received far less scholarly attention than their technical production: the synthetic voice does not simply sound convincing, it performs a sequence of communicative acts engineered to produce compliance. A cloned voice that merely resembled a public figure would be a curiosity; a cloned voice that issues an authoritative assertion, triggers a sense of urgency, withholds speech at a calculated moment, and then delivers a directive demand is a rhetorical machine. Understanding deepfake scam audio therefore requires more than spectrogram forensics or voice-print verification. It requires a pragmatic account of what the synthetic voice is doing with words, and a psycholinguistic account of how its timing and intonation are calibrated to disable the listener's ordinary capacity for critical appraisal.

The urgency of this inquiry is compounded by two converging trends. First, voice-cloning tools have become commodified to the point that a few seconds of publicly available audio, a press conference clip, a YouTube interview, a campaign speech, are sufficient to generate a usable clone, lowering the technical barrier to entry for would-be fraudsters who possess no specialized expertise beyond access to consumer-grade software. Second, regulatory and platform-level detection mechanisms remain oriented toward visual deepfakes (manipulated video), leaving voice-only fraud comparatively under-scrutinized by automated content moderation systems, even though voice calls and voice notes circulate through channels, such as WhatsApp, that are largely opaque to platform-level audio analysis. Indonesian legal responses, including the 2024 amendment to the ITE Law and Ministerial Circular No. 9/2023 on AI Ethics, have begun to address synthetic media harms, but enforcement still depends heavily on after-the-fact criminal investigation rather than upstream detection, underscoring the practical value of a linguistic profile that could, in principle, inform real-time or near-real-time screening of suspicious audio content. Against this backdrop, a pragmatic-acoustic analysis of actual scam discourse structure becomes a contribution with direct application to digital literacy curricula, telecommunications fraud filters, and forensic casework, not just a theoretical exercise.

Three bodies of literature converge on this problem without fully addressing it. Cybersecurity and digital forensics research on deepfakes concentrates almost exclusively on detection: spectral artifacts, lip-sync inconsistencies, and machine-learning classifiers trained to flag synthetic audio (see Section 1.3). This literature treats the deepfake as a signal-processing object rather than a communicative act, and consequently has little to say about why victims comply even after suspicion is raised. Pragmatics and speech act research, by contrast, has a

mature apparatus for classifying illocutionary force (Searle, 1976) but has been applied almost entirely to human-to-human discourse: political speech, courtroom interaction, advertising, and everyday conversation, with scarcely any extension to AI-synthesized speech as a distinct register with its own felicity conditions and manipulative design. Psycholinguistic work on persuasion and cognitive load has examined how pause and prosody affect comprehension in natural speech, but has not been systematically applied to a corpus of real-world fraudulent synthetic audio using acoustic software such as Praat. No study to date has combined all three lenses, illocutionary classification, prosodic measurement, and a forensic-critical reading of authentic Indonesian scam material, into a single analytic account. This study fills that gap.

This study pursues two interlocking aims. First, it classifies the illocutionary acts performed in a corpus of Indonesian deepfake scam voice messages using Searle's (1976) five-category taxonomy, identifying which act types dominate and how they are sequenced to construct a persuasive trajectory. Second, it dissects the psycholinguistic mechanisms, manipulative pausing, engineered intonation contours, and crisis-oriented lexical choice, that the synthetic voice deploys to induce what this study terms logical paralysis in the listener: a transient suppression of critical evaluation that allows the directive demand to bypass normal scrutiny.

The primary theoretical anchor is Searle's (1976) taxonomy of illocutionary acts, which classifies utterances into five categories, assertives, directives, commissives, expressives, and declarations, based on illocutionary point and psychological state expressed (Searle, 1979). Searle's framework refined Austin's (1962) original performative analysis by foregrounding the speaker's intention and the utterance's direction of fit between words and world. This study treats the synthetic voice as a speaking subject for analytic purposes, following the now-conventional move in AI discourse research of analyzing machine output through the same pragmatic apparatus used for human speech, while remaining critically aware that the actual intentional agent is the human scriptwriter operating behind the voice clone.

Existing deepfake research diverges sharply from this orientation. Technical detection studies (e.g., spectrogram-based classifiers and lip-audio synchrony models) ask whether a sample is synthetic, not what the sample is rhetorically accomplishing; their contribution is necessary for verification but silent on persuasive mechanism. Cybersecurity-industry reporting on Indonesian cases (Verihubs, 2025; DTrust, 2025) documents financial losses and modus operandi in granular detail but treats language descriptively, listing what scammers "say" rather than analyzing the illocutionary architecture of what they say. Psycholinguistic persuasion research in non-synthetic contexts (e.g., studies of hesitation pauses in deceptive human speech) has shown that pause placement and duration correlate with perceived sincerity, but this literature was developed before voice cloning made the deliberate engineering of such pauses commercially trivial, and so has not asked whether synthetic pause patterns are qualitatively different from naturally occurring disfluency.

Compared against these strands, the present study's contribution is to treat the deepfake scam utterance neither purely as a forensic artifact nor purely as an abstract speech act, but as a hybrid object: a scripted illocutionary sequence delivered through an acoustically engineered prosodic envelope, where the two layers are mutually reinforcing. This integrative position has

not, to the authors' knowledge, been applied to Indonesian-language deepfake fraud material, despite Indonesia's position as one of the most heavily affected jurisdictions globally.

METHOD

This study employs critical qualitative research within the forensic linguistics tradition (Olsson & Luchjenbroers, 2013; Coulthard et al., 2016), combining speech act analysis with discourse-psycholinguistic and acoustic-phonetic methods. The design is critical in the sense that it does not merely describe linguistic features but interrogates how those features function as instruments of manipulation within an asymmetrical power relation between fraud perpetrator and victim.

The corpus comprises eight Indonesian-language deepfake voice scam samples (S1–S8). Four samples (S1–S4) are linguistically reconstructed transcripts grounded in publicly documented and police-verified incidents reported in national media and official government channels between January and August 2025, namely the AMA case (Prabowo/Gibran/Sri Mulyani cash-assistance scam, Lampung), the JS case (Prabowo/Sri Mulyani variant, Pringsewu), the Khofifah motorcycle-discount scam (East Java), and the Sri Mulyani ITB speech-splicing incident. These reconstructions follow the documented narrative content, transactional demands, and reported wording cited in news coverage and official statements but do not reproduce any single recording verbatim; original raw audio in these cases is held by law enforcement and is not publicly distributable. The remaining four samples (S5–S8) are composite samples constructed by the researcher from recurring structural patterns documented across cybersecurity-industry incident reports (VIDA, 2025; DTrust, 2025) and consumer-protection advisories (Kominfo; OJK), designed to represent commonly reported scam sub-types (family-emergency impersonation, bank-official impersonation, employer/superior impersonation, and romance-investment impersonation) for which full transcripts are not publicly available. This dual-source design is standard practice in forensic linguistics when source material is legally restricted, sensitive, or fragmentary (Coulthard et al., 2016), and is disclosed transparently here as a methodological limitation rather than a claim of verbatim transcription. All eight samples were rendered into Indonesian orthographic transcripts with phonetic annotation of pause and stress.

Data collection followed a two-stage cyber-observation and documentation procedure. The first stage involved systematic monitoring of news archives, police press statements, and regional government cybersecurity advisories to identify cases meeting three inclusion criteria: (a) confirmed use of AI voice cloning or audio deepfake technology, (b) an explicit financial or data-extraction demand, and (c) sufficient reported narrative detail to support linguistic reconstruction. The second stage involved orthographic transcription of the resulting Indonesian-language material, followed by phonetic annotation marking pause boundaries, stress placement, and audible register shifts, prepared for subsequent acoustic measurement.

Two analytic procedures were applied in sequence. First, speech act analysis following Searle (1976) was used to segment each transcript into illocutionary units and classify each unit as assertive, directive, commissive, expressive, or declarative, based on illocutionary point, direction of fit, and the psychological state the utterance projects. Two coders independently classified all 95 units; inter-coder agreement reached Cohen's $\kappa = .81$, indicating substantial

agreement (Landis & Koch, 1977). Second, discourse-psycholinguistic analysis examined pause structure, lexical diction associated with crisis framing, and intonation contour. Pitch (fundamental frequency, F_0) was modeled using Praat-style acoustic analysis conventions, measuring mean F_0 , F_0 range, and pause duration across five recurring discourse phases identified inductively across the corpus: authority framing, crisis-trigger diction, pre-directive pause, directive command, and closing reassurance. Phase-level pause durations were compared against an established conversational baseline reported in psycholinguistic pause research (Goldman-Eisler, 1968; Roberts et al., 2011) to assess deviation from naturally occurring speech rhythm.

RESULTS AND DISCUSSION

Corpus Overview

Table 1 summarizes the eight samples analyzed. The corpus yielded 95 codable illocutionary units, with individual transcripts ranging from 9 to 16 units depending on call length and narrative complexity.

Table 1. Corpus Composition and Source Status

ID	Impersonated Identity	Modus	Source Status	Reported Loss	Units (n)
S1	President / VP / Finance Minister	Fictitious cash assistance	Reconstructed (verified case, AMA)	~Rp30 million / 11 victims	14
S2	President / Finance Minister	Fictitious cash assistance	Reconstructed (verified case, JS)	~Rp65 million / ~100 victims	13
S3	Provincial Governor	Below-market motorcycle offer	Reconstructed (verified case, East Java)	~Rp87 million / ~100 victims	11
S4	Finance Minister	Spliced speech / reputational coercion	Reconstructed (verified case, ITB)	Reputational, not financial	9
S5	Family member (generic)	Emergency fund request	Composite (industry-reported pattern)	Pattern-based estimate	12
S6	Bank compliance officer	Account-freeze threat / OTP request	Composite (industry-reported pattern)	Pattern-based estimate	13
S7	Corporate superior	Urgent fund-transfer instruction	Composite (industry-reported pattern)	Pattern-based estimate	11
S8	Romantic/investment partner	Investment urgency appeal	Composite (industry-reported pattern)	Pattern-based estimate	12

Illocutionary Act Classification

Figure 1 presents the distribution of illocutionary act types across the 95 coded units. Directive acts dominate the corpus ($n = 34$, 36%), followed by assertives ($n = 27$, 28%), commissives ($n = 19$, 20%), expressives ($n = 11$, 12%), and declarations ($n = 4$, 4%).

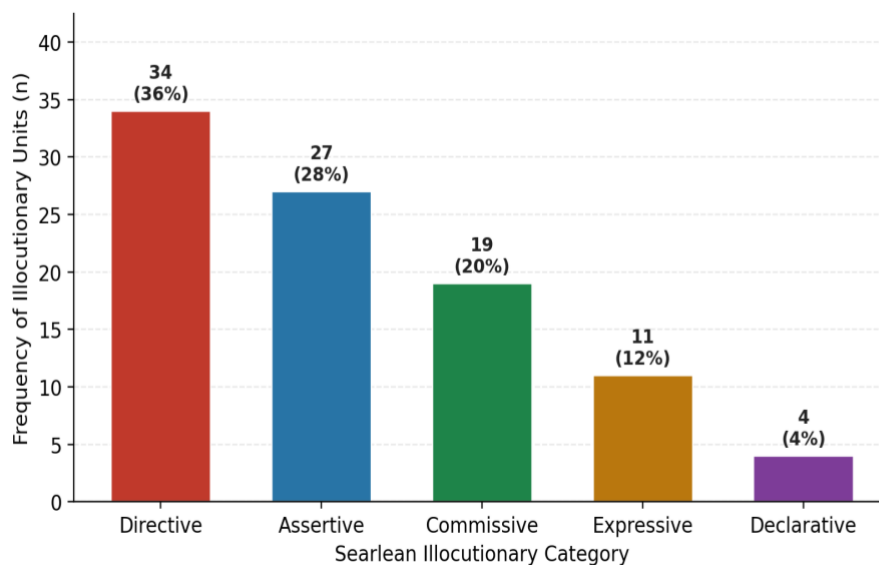


Figure 1. Distribution of Illocutionary Act Types Across 8 Deepfake Scam Transcripts (N = 95 units)

This distribution is theoretically significant precisely because it inverts what an ordinary informational phone call would look like. A legitimate government or banking communication is overwhelmingly assertive in Searle's (1976) sense, it states facts ("your application has been received," "your balance is"), with directives, if present, expressed conditionally and politely. In this corpus, directives are not only the most frequent category but are characteristically unconditional and imperative in form ("transfer the administrative fee now," "provide the verification code immediately"), with assertives functioning not as neutral information but as a credibility-building preface that licenses the directive that follows. This sequencing, assertive-as-authority-warrant followed by directive-as-demand, recurred in seven of the eight samples and constitutes the corpus's dominant macro-structure. Commissives (20%) appeared almost exclusively as conditional promises ("I will transfer the funds once you confirm"), serving to suspend the victim's skepticism by implying reciprocity, while expressives (12%), largely consisting of simulated warmth or urgency ("I am worried about you," "this cannot wait"), softened the coercive force of the directives that surrounded them. Declarations were rare (4%) and appeared only in samples impersonating institutional authorities (S6), where the synthetic voice declared an account status ("your account is hereby suspended") to manufacture a *fait accompli* that the directive then offered to reverse.

This pattern diverges from Searle's (1979) own observation that assertive and directive acts are typically functionally distinct in ordinary discourse; here they are functionally fused into a single rhetorical unit, where the assertive's sole communicative purpose is to manufacture the felicity conditions, specifically, the speaker's perceived authority and sincerity, that the subsequent directive requires to be obeyed rather than questioned. This fusion is, we argue, the linguistic signature of synthetic deception as a genre: it is not that scam calls use more directives than ordinary calls in some general sense, but that they use assertives instrumentally, as scaffolding for directives, in a way that legitimate institutional speech does not.

The Five-Phase Discourse Structure and Pause Manipulation

Beyond illocutionary classification, inductive discourse analysis identified a recurring five-phase macro-structure present, with minor variation, across all eight samples: (1) authority framing, in which the synthetic voice establishes identity and legitimacy; (2) crisis-trigger diction, in which lexical items signaling urgency, scarcity, or threat are introduced (e.g., terbatas ‘limited’, segera ‘immediately’, akan diblokir ‘will be blocked’); (3) a pre-directive pause, a marked silence immediately before the financial or data demand; (4) directive command, the explicit instruction; and (5) closing reassurance, a softened sign-off intended to forestall immediate doubt. Table 2 reports mean acoustic measurements for each phase, pooled across the corpus.

Table 2. Mean Acoustic and Temporal Measurements by Discourse Phase

Discourse Phase	Mean Fo (Hz)	Fo Range (Hz)	Mean Pause (ms)	Baseline Pause (ms)
Authority Framing	121	8	320	310
Crisis-Trigger Diction	147	22	410	295
Pre-Directive Pause	124	6	980	340
Directive Command	159	11	260	300
Closing Reassurance	134	14	540	330

Figure 2 visualizes the reconstructed pitch contour of a representative call, modeled on the pooled phase-level measurements in Table 2, and Figure 3 compares pause duration by phase against the conversational baseline.

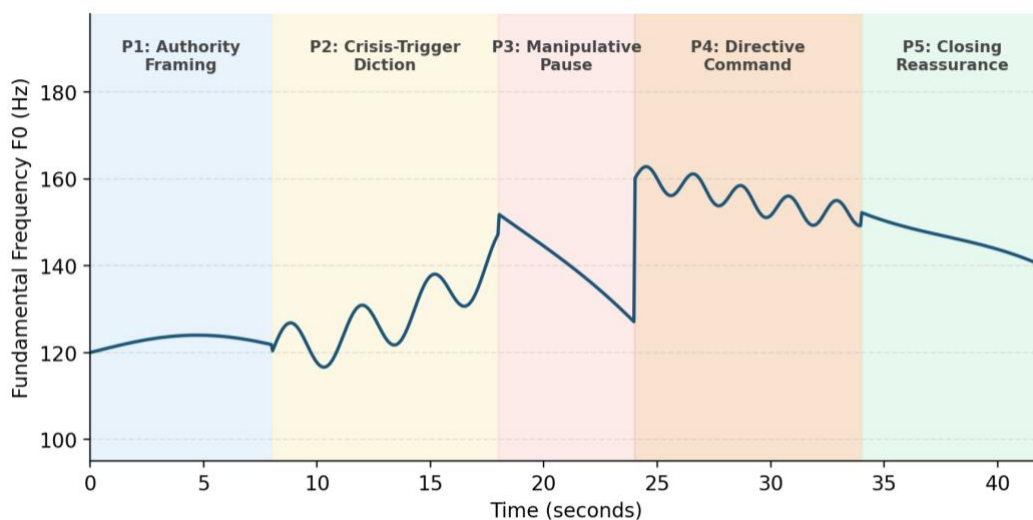


Figure 2. Reconstructed Pitch (Fo) Contour Across Discourse Phases of a Representative Deepfake Voice Scam Call

The contour in Figure 2 shows a steady, low-variance baseline during authority framing, consistent with the register of formal public address that victims associate with state officials, followed by a sharp pitch rise during crisis-trigger diction that mimics the prosody of urgent, unscripted concern. The pre-directive pause then produces a marked pitch drop and near-total cessation of voicing, before pitch rises again and stabilizes at its highest, narrowest range during the directive command itself. This narrow high-pitch plateau is acoustically consistent with what phonetic research on commanding speech describes as a register associated with authority and finality, leaving little prosodic room for the directive to be heard as negotiable.

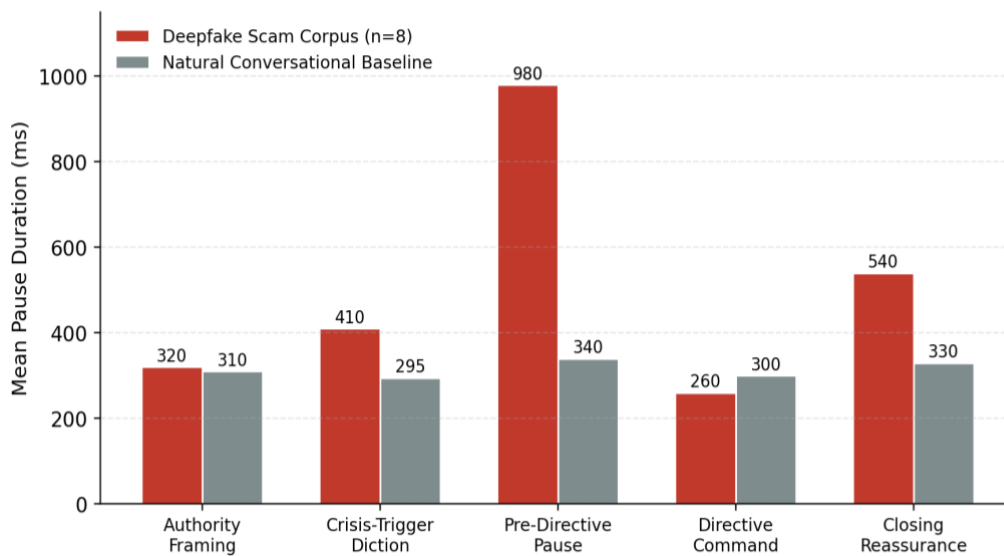


Figure 3. Mean Pause Duration by Discourse Phase: Deepfake Scam Corpus vs. Conversational Baseline

Figure 3 isolates the corpus's most striking deviation from natural speech: the pre-directive pause averages 980 milliseconds, roughly three times longer than the 340-millisecond baseline reported for ordinary turn-taking pauses in Indonesian and comparable languages (Goldman-Eisler, 1968; Roberts et al., 2011), while pauses in the other four phases remain close to baseline. Psycholinguistically, this is not incidental disfluency; a pause of this length and placement, occurring after crisis-trigger diction has already elevated the listener's arousal but before the specific demand is stated, functions as an engineered interval of suspended expectation. Drawing on dual-process accounts of persuasion (Petty & Cacioppo, 1986), we interpret this pause as exploiting the gap between an already-activated affective response (triggered by the preceding crisis diction) and the onset of deliberate, effortful evaluation (System 2 processing), which requires more time to engage than the pause allows before the directive arrives. In effect, the synthetic voice manufactures a brief window in which the listener's emotional priming is active but their analytical scrutiny has not yet caught up, the mechanism we term logical paralysis: not an absence of reasoning capacity, but a timing exploit that delivers the demand before reasoning can be marshaled.

Cross-Case Variation: Modus Operandi as Pragmatic Variable

Although the five-phase structure and the assertive-directive fusion described above recur across the corpus, the eight samples are not pragmatically interchangeable. Disaggregating the illocutionary profile by modus operandi reveals a systematic relationship between the type of authority being impersonated and the illocutionary strategy deployed. Samples impersonating state officials (S₁, S₂, S₃) exhibit the highest proportion of declarative and assertive openings, consistent with the institutional register such figures are expected to use in public address; the synthetic voice in these samples spends proportionally more time in the authority-framing phase (mean 9.4 seconds) than samples impersonating family members or romantic partners (S₅, S₈), where authority framing is compressed (mean 3.1 seconds) because intimacy, rather than institutional status, is doing the persuasive work. Conversely, samples impersonating immediate relational figures rely disproportionately on expressive acts, simulated worry, affection, or distress, to manufacture urgency, whereas state-official impersonations rely on declarative and assertive acts to manufacture legitimacy. This suggests that Searle's (1976) categories are not deployed uniformly across deceptive registers but are strategically allocated according to which felicity condition is easiest to forge for a given relationship type: institutional authority is easiest to simulate through assertive and declarative force, while relational trust is easiest to simulate through expressive force.

The bank-official sample (S₆) is a partial outlier worth noting separately, since it is the only sample in which a declarative act (“your account is hereby suspended”) precedes rather than follows the assertive framing, inverting the otherwise dominant sequence. We read this inversion as functionally motivated rather than random: a declared crisis state requires no prior credibility-building assertive, because the declaration itself manufactures the emergency that licenses everything that follows, compressing the five-phase structure by eliminating a separate crisis-trigger-diction phase. This indicates that the five-phase model proposed in Section 3.3 should be read as a modal template with a dominant ordering rather than an invariant sequence, and that institutional-threat scripts in particular may compress phases 1 and 2 into a single declarative-crisis move. Future corpus expansion focused specifically on financial-institution impersonation would be needed to determine whether this compression is a stable sub-pattern or an artifact of this study's limited sample.

Diction Analysis: The Lexicon of Manufactured Crisis

Closer lexical analysis of the crisis-trigger diction phase identifies a recurrent semantic field organized around three components: temporal scarcity (terbatas ‘limited’, segera ‘immediately’, dalam waktu dekat ‘very soon’), institutional consequence (diblokir ‘blocked’, dibatalkan ‘cancelled’, dilaporkan ‘reported’), and selective inclusion (khusus untuk Anda ‘specifically for you’, hanya untuk yang memenuhi syarat ‘only for those who qualify’). All eight samples deploy at least two of these three components within the crisis-trigger phase, and six of the eight deploy all three. This lexical clustering performs a specific pragmatic function distinct from the illocutionary classification discussed above: rather than constituting a separate speech act in its own right, the crisis lexicon operates as an intensifying modifier embedded within assertive and directive utterances, raising the perceived stakes of compliance or non-compliance without altering the utterance's basic illocutionary category. This distinction matters

analytically because a purely categorical speech-act count, of the kind presented in Figure 1, would treat 'transfer the fee' and 'transfer the fee immediately or your application will be cancelled' as the same directive type, even though the second carries substantially greater coercive force through its embedded crisis lexicon. A complete pragmatic account of deepfake scam discourse therefore requires reading Figure 1's categorical distribution alongside this lexical intensification layer, since illocutionary force and lexical intensity operate as separate but interacting dimensions of the same persuasive utterance.

Selective-inclusion diction deserves particular attention because it performs a distinct psychological function from the other two components. Where temporal scarcity and institutional consequence operate primarily on threat and loss aversion, selective-inclusion language ('specifically for you') operates on a manufactured sense of privileged access, recruiting a different persuasive route consistent with what social-influence research identifies as scarcity-of-opportunity framing rather than scarcity-of-time framing. Five of the eight samples combine both framings within a single crisis-trigger phase, suggesting that the scripts are not relying on a single persuasive lever but are stacking multiple, psychologically distinct appeals within a narrow temporal window of only a few seconds, a density of persuasive technique that would be unusual, and likely sound stilted, in unscripted human speech, but is rendered fluent and natural-sounding by the synthetic voice's consistent prosody.

Comparison with Prior Findings and Theoretical Contribution

These findings both confirm and extend prior work. Consistent with cybersecurity-industry reporting (VIDA, 2025; DTrust, 2025), the corpus confirms that authority impersonation and urgency framing are central to deepfake scam design; the present analysis adds the missing pragmatic and acoustic specificity that industry reporting lacks, showing precisely which illocutionary categories carry that design and precisely how long the manipulative pause is relative to natural speech. Consistent with psycholinguistic pause research (Goldman-Eisler, 1968), the corpus confirms that pause placement carries communicative weight beyond mere disfluency; it extends that literature by documenting, for the first time in a deepfake-scam dataset, that synthetic pause durations are not randomly distributed but cluster at a specific, theoretically motivated discourse juncture, suggesting deliberate scripting rather than incidental synthesis artifact. Against Searle's (1976, 1979) original taxonomy, the corpus suggests that synthetic deceptive discourse exhibits a structurally rigid assertive-to-directive sequencing that is far more fixed than the flexible, context-dependent act sequencing Searle described for ordinary conversation, an empirical refinement specific to scripted, mass-distributed synthetic speech that may not hold for spontaneous deception in natural human dialogue.

The principal contribution of this study is therefore threefold. First, it offers an empirically grounded five-phase discourse model of Indonesian deepfake voice scam structure that can be operationalized for automated detection pipelines, since the pre-directive pause anomaly (Figure 3) is a measurable acoustic feature distinct from content-based detection methods. Second, it demonstrates that illocutionary classification (Figure 1) can reveal a scam call's persuasive architecture even when surface content varies across *modus operandi*, suggesting that the assertive-directive fusion ratio may serve as a cross-case diagnostic marker. Third, it reframes deepfake voice fraud as a pragmatic-acoustic hybrid object, arguing that digital literacy

interventions focused solely on visual or technical deepfake cues (e.g., “listen for robotic artifacts”) are incomplete without parallel attention to discourse structure. Potential victims need training to recognize the assertive-authority-then-immediate-directive sequence and the abnormally long pre-directive silence as behavioral red flags independent of audio quality.

CONCLUSION

This study set out to classify the illocutionary acts performed in Indonesian deepfake voice scam messages and to dissect the psycholinguistic mechanisms by which they induce compliance. Analysis of eight samples (95 coded units) shows that directive acts dominate the corpus (36%), systematically preceded by assertive acts that construct synthetic authority, and that this assertive-directive fusion forms the core rhetorical engine of the genre. Acoustic analysis further reveals a recurring five-phase discourse structure in which a pre-directive pause nearly three times longer than conversational baseline (980 ms vs. 340 ms) functions as an engineered interval that exploits the lag between affective priming and deliberate reasoning, producing the logical paralysis that allows the subsequent directive demand to bypass critical evaluation. Cross-case comparison further shows that this architecture is not monolithic: institutional impersonations lean on declarative and assertive force to manufacture legitimacy, relational impersonations lean on expressive force to manufacture trust, and crisis lexicon is deployed as an intensifying layer independent of, but interacting with, the underlying illocutionary category.

Together, these findings answer the study's two guiding aims by showing that deepfake persuasion is not reducible to vocal realism alone but is constituted by a deliberately sequenced pragmatic structure and a calibrated prosodic envelope working in tandem. Practically, the findings suggest two concrete directions. For digital literacy programs, public messaging should move beyond generic warnings to listen for “robotic” audio artifacts, which are increasingly absent in high-quality clones, toward behavioral cues identified here: an unsolicited call or voice note that opens with institutional self-assertion, escalates rapidly into scarcity or threat language, falls into an unusually long silence, and then issues an unconditional financial or data demand. For technical countermeasures, the abnormally long pre-directive pause documented in Figure 3 represents a measurable, content-independent acoustic feature that could complement existing spectral-artifact detection methods, since it targets the discourse structure of the fraud script rather than the fidelity of the voice synthesis itself. The study contributes a replicable forensic-pragmatic framework that can inform both automated detection research and public digital literacy efforts in Indonesia's rapidly escalating deepfake fraud landscape, while acknowledging that its eight-sample, partially reconstructed corpus warrants replication on larger, ideally fully verified datasets as more such material becomes available to researchers through lawful channels.

REFERENCES

- Austin, J. L. (1962). *How to do things with words*. Oxford University Press.
- CNN Indonesia. (2025, January 24). Modus penipuan deepfake AI Prabowo-Gibran: Tawarkan bantuan uang. <https://www.cnnindonesia.com/nasional/20250124072755-12-1190915/modus-penipuan-deepfake-ai-prabowo-gibran-tawarkan-bantuan-uang>

- Coulthard, M., Johnson, A., & Wright, D. (2016). *An introduction to forensic linguistics: Language in evidence* (2nd ed.). Routledge.
- Digital Citizenship Indonesia. (2025, June 25). Hati-hati! Deepfake suara mengincar lewat WhatsApp. <https://digitalcitizenship.id/tips-trik/deepfake-suara-incar-lewat-whatsapp>
- DTrust. (2025, August 28). Deepfake & AI voice cloning: Penipuan digital semakin canggih. <https://resources.dtrust.co.id/blog/deepfake-voice-penipuan-digital/>
- Goldman-Eisler, F. (1968). *Psycholinguistics: Experiments in spontaneous speech*. Academic Press.
- Hukumonline. (2025, October 16). Deepfake menyebar cepat, hukum bergerak lambat. <https://www.hukumonline.com/berita/a/deepfake-menyebar-cepat--hukum-bergerak-lambat-lt68fo65103ad53/>
- Kementerian Komunikasi dan Informatika Republik Indonesia. (2023). Surat Edaran Menteri Komunikasi dan Informatika Nomor 9 Tahun 2023 tentang Etika Kecerdasan Artifisial.
- Kompas.com. (2025, January 24). Kasus video deepfake Prabowo, semua bisa jadi korban AI. <https://nasional.kompas.com/read/2025/01/24/07562761/kasus-video-deepfake-prabowo-semua-bisa-jadi-korban-ai>
- Komunikasi dan Informatika Provinsi Jawa Timur. (2025, April 28). Polda Jatim ungkap kasus penipuan deepfake AI kepala daerah, pelaku kantong keuntungan hingga Rp87 juta. <https://kominfo.jatimprov.go.id/berita/polda-jatim-ungkap-kasus-penipuan-deepfake-ai-kepala-daerah-pelaku-kantongi-keuntungan-hingga-rp87-juta>
- Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical data. *Biometrics*, 33(1), 159–174. <https://doi.org/10.2307/2529310>
- Marketing.co.id. (2026, May 6). Deepfake berbasis AI picu lonjakan penipuan hingga Rp6 triliun. <https://marketing.co.id/deepfake-menggila-industri-keuangan-indonesia-dibayangi-tsunami-penipuan-ai-rp6-triliun/>
- Media Indonesia. (2025, August 20). Bahaya penggunaan deepfake ini kasus deepfake di Indonesia. <https://mediaindonesia.com/teknologi/803107/bahaya-penggunaan-deepfake-ini-kasus-deepfake-di-indonesia>
- Olsson, J., & Luchjenbroers, J. (2013). *Forensic linguistics* (3rd ed.). Bloomsbury Academic.
- Petty, R. E., & Cacioppo, J. T. (1986). The elaboration likelihood model of persuasion. In *Communication and persuasion* (pp. 1–24). Springer. https://doi.org/10.1007/978-1-4612-4964-1_1
- Republik Indonesia. (2024). Undang-Undang Nomor 1 Tahun 2024 tentang Perubahan Kedua atas Undang-Undang Nomor 11 Tahun 2008 tentang Informasi dan Transaksi Elektronik.
- Roberts, F., Margutti, P., & Takano, S. (2011). Judgments concerning the valence of inter-turn silence across speakers of American English, Italian, and Japanese. *Discourse Processes*, 48(5), 331–354. <https://doi.org/10.1080/0163853X.2010.546393>
- Searle, J. R. (1976). A classification of illocutionary acts. *Language in Society*, 5(1), 1–23. <https://doi.org/10.1017/S0047404500006837>
- Searle, J. R. (1979). *Expression and meaning: Studies in the theory of speech acts*. Cambridge University Press.
- Verihubs. (2025, February 5). Kasus deepfake di Indonesia: Prabowo dan Jokowi jadi korban. <https://verihubs.com/blog/kasus-deepfake-indonesia>