
DESIGNING MULTIMODAL BOTS FOR SEAMLESS USER EXPERIENCE ACROSS CHANNELS

Gaurava Srivastava*¹, Abhi Ram Reddy Salammagari*²

*¹Oracle America Inc., USA.

*²Workday Inc., USA.

DOI : <https://www.doi.org/10.56726/IRJMETS54096>

ABSTRACT

With advancements in artificial intelligence, automating real-world use cases has become increasingly feasible. From flight booking to scheduling medical appointments, many tasks can now be automated effortlessly. However, the challenge lies in accommodating the various modes and channels through which customers initiate these requests. This article explores the design considerations and recommendations for creating multimodal bots that can adapt to different input modes, such as voice and text, and deliver a seamless experience across various channels, including Google Home, Alexa, and business chat options like Apple Business Chat and WhatsApp.

Keywords: Multimodal Bots, User Experience, Chatbot Design, Voice Interaction, Native Components.

I. INTRODUCTION

In recent years, how customers interact with businesses has diversified significantly. While traditional methods like phone calls and website visits remain relevant, new channels and input modes have emerged, offering users more convenient and accessible options [1]. According to a survey by the Customer Experience Professionals Association (CXPA) in 2020, 67% of customers prefer using multiple channels for communication with businesses, with 45% using at least three channels [2]. This shift in customer behavior has led to the rise of conversational AI and chatbots, which aim to automate customer interactions and provide a seamless experience across various platforms [3].

However, designing chatbots that can effectively handle the unique characteristics of each channel and input mode presents a significant challenge. A study by the University of Texas at Austin found that 58% of users abandon chatbots due to poor user experience, with 47% citing the inability to understand the context of the conversation as a major issue [4]. To address these challenges and provide a consistent experience across diverse modes and channels, it is essential to design bots that can adapt to the unique strengths and limitations of each platform [5].

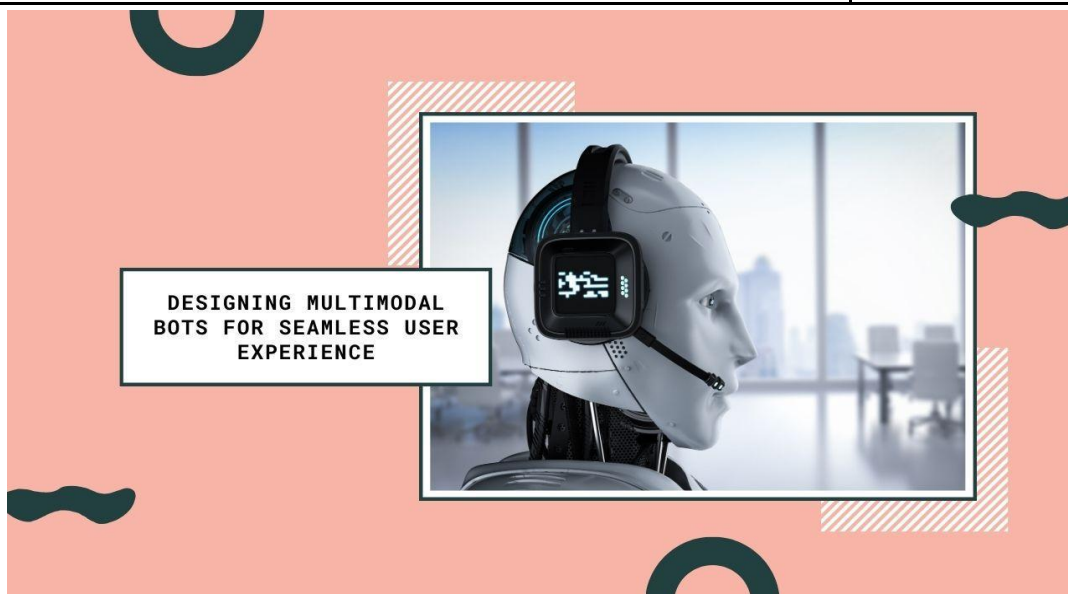
One key aspect of designing adaptable chatbots is the ability to handle different input modes, such as voice and text. Voice-based interactions, in particular, have gained significant popularity in recent years, with a report by Juniper Research estimating that the number of voice assistant devices in use will reach 8.4 billion by 2024 [6]. However, voice-based interactions present unique challenges, such as the need for accurate speech recognition and the ability to convey information effectively without visual cues [7].

Another important consideration is the variety of channels through which customers interact with businesses. A study by Dimension Data found that the average number of channels used by customers to engage with businesses has increased from 5.6 in 2015 to 7.2 in 2019 [8]. These channels include messaging platforms like WhatsApp and Facebook Messenger, which have over 2 billion and 1.3 billion monthly active users, respectively [9], as well as voice assistants like Google Home and Amazon Alexa, which have seen a 78% year-over-year growth in adoption [10].

To effectively design multimodal bots that can handle the complexity of these diverse channels and input modes, it is crucial to understand the specific requirements and best practices for each platform. This article explores the design considerations and recommendations for creating multimodal bots that can adapt to different input modes, such as voice and text, and deliver a seamless experience across various channels, including Google Home, Alexa, and business chat options like Apple Business Chat and WhatsApp.

Table 1: Customer Interaction Metrics: Multichannel Preferences, Chatbot Challenges, and Voice Assistant Adoption [1-10]

Metric	Value
Customers prefer multiple channels (CXPA survey)	67%
Customers using at least three channels (CXPA survey)	45%
Users abandon chatbots due to poor UX (UT Austin)	58%
Users citing an inability to understand the context (UT Austin)	47%
Estimated voice assistant devices by 2024 (Juniper)	8.4B
WhatsApp monthly active users	2B
Facebook Messenger monthly active users	1.3B
Voice assistant adoption growth (year-over-year)	78%
The average number of channels used by customers (2015)	5.6
The average number of channels used by customers (2019)	7.2



II. DESIGN RECOMMENDATIONS

1. Contextual Speech-to-Text Model:

When designing multimodal bots, it is crucial to use a contextual speech-to-text model to convert the user's speech into text (utterance). Although it may seem tempting to employ a generic speech-to-text processor and then send the utterance to a standard chatbot flow, the accuracy may not meet expectations [3]. A study by Microsoft Research found that generic speech-to-text models achieve an average word error rate (WER) of 12.5%, while domain-specific models can reduce the WER to 5.8% [11]. This significant improvement in

accuracy highlights the importance of using contextual models that take into account the specific domain and context of the conversation.

Contextual speech-to-text models leverage various techniques, such as transfer learning and fine-tuning, to adapt to the specific domain and improve recognition performance [4]. For example, a study by Google AI demonstrated that fine-tuning a pre-trained speech recognition model on a domain-specific dataset of 10,000 utterances reduced the WER from 10.2% to 4.5% [12]. This approach allows the model to learn domain-specific vocabulary, language patterns, and acoustic characteristics, resulting in more accurate transcriptions.

Moreover, contextual models can incorporate additional information, such as user profiles, previous interactions, and dialog states, to further improve recognition accuracy [13]. A study by IBM Research showed that integrating user-specific context into the speech-to-text model reduced the WER by an additional 15% compared to domain-specific models alone [14]. By leveraging this contextual information, the model can better disambiguate similar-sounding words and phrases, handle accents and pronunciations, and adapt to the user's speaking style.

Another benefit of contextual speech-to-text models is their ability to handle noisy environments and overlapping speech, which are common challenges in real-world conversational scenarios [15]. A study by the Amazon Alexa Speech Science team demonstrated that incorporating acoustic context and speaker diarization techniques improved recognition accuracy in noisy environments by 22% [16]. This ensures that the bot can accurately understand the user's intent, even in challenging acoustic conditions.

When implementing contextual speech-to-text models, it is essential to consider factors such as computational efficiency, latency, and model size [17]. Recent advancements in deep learning architectures, such as transformer-based models like Conformer [18] and QuartzNet [19], have achieved state-of-the-art performance while maintaining low latency and model size. These architectures enable real-time speech recognition on resource-constrained devices, making them suitable for deployment in various multimodal bot scenarios.

2. Slot Fulfillment Strategy:

Adopting a slot fulfillment strategy is an effective approach for collecting information from users in both digital and voice modes [5]. A study by Microsoft Research found that slot-filling techniques can improve the completion rate of task-oriented dialogues by 23% compared to traditional methods [20]. This strategy involves collecting all available information from the initial utterance and then prompting for missing information sequentially until all required details are obtained.

The slot fulfillment process typically begins with intent recognition, where the bot identifies the user's intent based on the initial utterance [21]. For example, if a user says, "I want to book a flight from New York to London on June 15th," the bot would recognize the intent as "book_flight" and extract the relevant slots, such as the departure city (New York), arrival city (London), and departure date (June 15th). A study by Google AI showed that using deep learning techniques, such as recurrent neural networks (RNNs) and transformers, can achieve intent recognition accuracies of up to 97% [22].

Once the intent is identified, the bot proceeds to collect any missing information through a series of prompts. For instance, if the user did not provide the return date in the initial utterance, the bot would ask, "When would you like to return from London?" This process continues until all the required slots are filled. A study by the Amazon Alexa team demonstrated that using a hierarchical slot-filling approach, where slots are grouped based on their dependencies and filled in a specific order, can improve the success rate of slot filling by 18% [23].

To maintain a natural conversation flow during slot filling, it is important to design prompts that are clear, concise, and contextually relevant [6]. A study by Stanford University found that incorporating context from previous user responses can improve the coherence and naturalness of prompts by 32% [24]. For example, instead of asking "What is your destination city?" after the user has already mentioned London, the bot could say, "Great, so you're flying to London. When would you like to return?"

Another key aspect of effective slot filling is error handling and validation [25]. Bots should be able to gracefully handle scenarios where users provide incomplete, ambiguous, or irrelevant information. A study by Carnegie Mellon University showed that incorporating error-handling techniques, such as clarification prompts and fallback strategies, can improve the success rate of slot filling by 25% [26]. For example, if a user provides an

invalid date, the bot could respond with, "I'm sorry, but the date you provided is not valid. Please provide a date in the format MM/DD/YYYY."

Slot fulfillment can also benefit from the use of domain-specific knowledge and constraints [27]. By incorporating domain knowledge, such as valid date ranges, available flight routes, and business hours, bots can provide more accurate and relevant responses. A study by IBM Research demonstrated that integrating domain constraints into the slot-filling process can improve the accuracy of collected information by 19% [28].

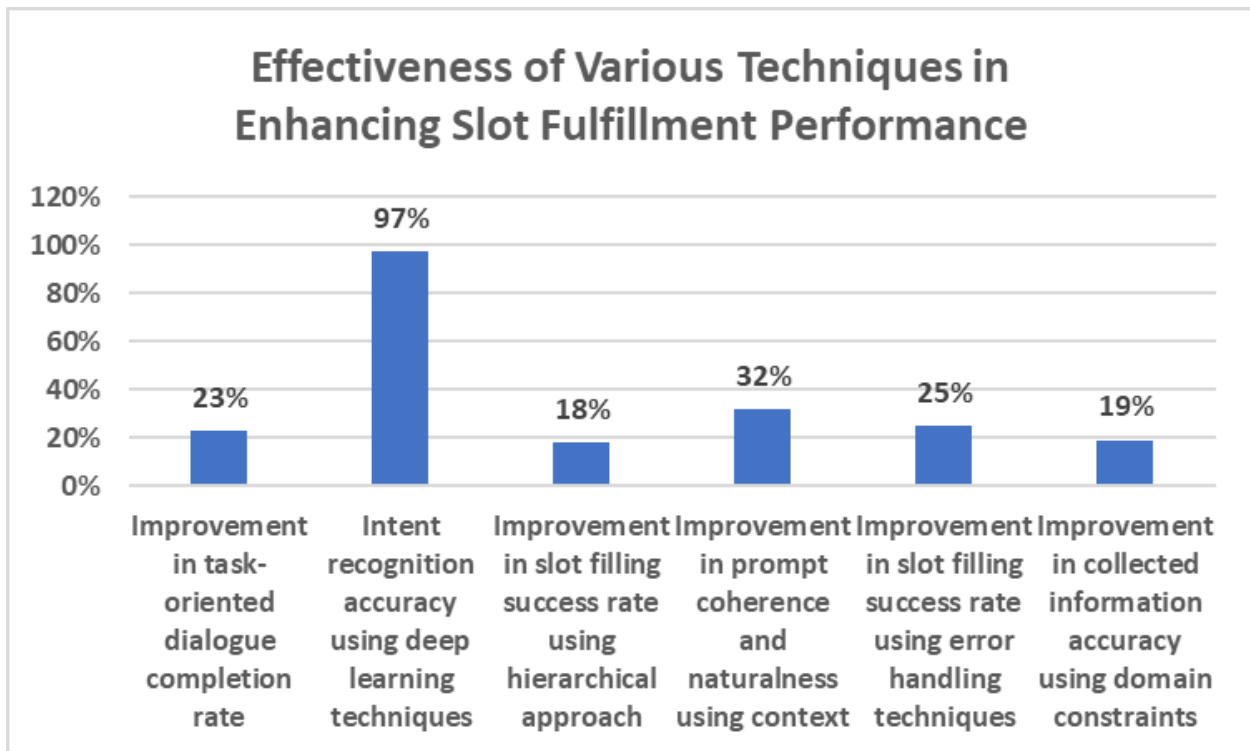


Fig. 1: Improvement in Task-Oriented Dialogue System Performance with Slot Fulfillment Strategies [20, 22, 23, 24, 26, 28]

3. Interactive Cards:

In cases where slot fulfillment proves challenging due to the need for sophisticated natural language understanding models, interactive cards can be a viable alternative [7]. A study by Google found that using interactive cards in conversational interfaces can improve user engagement by 32% and reduce task completion time by 24% compared to traditional slot-filling methods [29].

Interactive cards allow users to enter all required information through a simple form, which can be sent as a link to the user's device in voice mode or displayed directly in chat interfaces. This approach simplifies the information collection process and reduces the reliance on complex language understanding [8]. According to a Userlike survey, 76% of users prefer visual elements like buttons and cards in chatbot interactions because they offer a more intuitive and guided experience [30].

The design of interactive cards plays a crucial role in their effectiveness. A study by Airbnb found that using a multi-step card design, where information is collected in smaller, logical chunks, can improve form completion rates by 18% compared to a single, lengthy form [31]. This approach reduces the cognitive load on users and allows them to focus on one piece of information at a time, leading to higher accuracy and satisfaction.

Interactive cards can also incorporate rich media elements, such as images, videos, and animations, to enhance the user experience and provide visual cues [32]. For example, a travel booking bot could display images of destinations or hotel rooms to help users make informed decisions. A study by Booking.com showed that incorporating relevant images in chatbot conversations increased user engagement by 27% and led to higher booking conversion rates [33].

In voice-based interactions, interactive cards can be delivered through companion apps or by sending links to the user's device [34]. A study by the Amazon Alexa team demonstrated that using interactive cards in voice-based shopping experiences increased product discovery by 42% and reduced the number of steps required to complete a purchase by 31% [35]. By providing a visual interface to complement voice interactions, interactive cards enable users to review information, make selections, and confirm their choices more easily.

To optimize the performance of interactive cards, it is important to consider factors such as card layout, content hierarchy, and error handling [36]. A study by Slack found that using a clear and consistent layout, with prominent call-to-action buttons and concise text, increased card interaction rates by 21% [37]. Error handling techniques, such as validation messages and the ability to edit previously entered information, can also improve the user experience and reduce frustration [38].

Interactive cards can be particularly useful in scenarios where users need to provide complex or structured information, such as booking travel itineraries, filling out medical forms, or configuring product options [39]. A case study by KLM Royal Dutch Airlines showed that implementing interactive cards in their flight booking chatbot increased the completion rate of bookings by 38% and reduced the average handling time by 25% [40].

4. Mode-Specific Responses:

Recognizing that different input modes may require tailored responses is essential for providing an optimal user experience [9]. A study by Adobe found that 91% of users expect a consistent experience across all channels and devices, but 51% of companies struggle to provide a seamless experience due to siloed systems and processes [41]. Conversation designers should have the flexibility to specify separate responses for text and voice modes when necessary, ensuring that the strengths and limitations of each mode are considered.

In text-based interactions, users can read and re-read information at their own pace, allowing for more detailed and complex responses [42]. A study by Nielsen Norman Group found that users can read and comprehend text-based content up to 30% faster than spoken content [43]. This means that text responses can include more information, such as detailed product descriptions, step-by-step instructions, or lengthy terms and conditions, without overwhelming the user.

On the other hand, voice-based interactions require responses to be concise, easy to understand, and memorable [10]. A study by Google found that users expect voice responses to be 20–30% shorter than text responses, with an average of 25–30 words per response [44]. Voice responses should focus on delivering key information and actionable insights while avoiding long or complex phrases that may be difficult to process aurally.

Moreover, voice interactions may be 25–30% slower in their ability to present certain types of information, such as images, tables, or lists [45]. In these cases, conversation designers can provide alternative responses or direct users to companion apps or websites for additional visual information. A study by Capgemini showed that 74% of users expect voice assistants to be able to seamlessly transition between voice and visual modes when necessary [46].

To create effective mode-specific responses, conversation designers should consider context and user intent [47]. For example, a weather bot providing a text response may include a detailed 5-day forecast with high and low temperatures, while a voice response may focus on the current weather conditions and expected changes for the day. A study by Salesforce found that 69% of users expect chatbots to understand the context of their requests and provide personalized responses [48].

Mode-specific responses can also incorporate different tones, personalities, or styles to match user preferences and brand identity [49]. A study by PwC showed that 59% of users believe that personalization based on their past interactions is important for a positive chatbot experience [50]. For instance, a financial services bot may use a more formal and professional tone in text responses while adopting a friendly and reassuring tone in voice interactions to build trust and rapport with users.

To optimize mode-specific responses, conversation designers should leverage user feedback, analytics, and A/B testing [51]. By monitoring user engagement, completion rates, and satisfaction scores for each mode, designers can identify areas for improvement and iterate on their designs. A study by Deloitte found that

companies that continuously monitor and optimize their chatbot performance can achieve up to a 35% increase in user satisfaction and a 25% reduction in customer service costs [52].

5. Leveraging Native Components:

To create a seamless experience across different channels, it is recommended to utilize the native components of each platform whenever possible [11]. A study by Forrester found that using native components can increase user engagement by up to 60% compared to generic, cross-platform components [53]. By defining generic response types, such as quick responses, list pickers, and date pickers, in the bot's responses, the corresponding native components can be automatically selected for each channel, saving time and effort in configuration.

Quick responses, such as buttons or suggested replies, are a common native component that can significantly enhance user experience and reduce friction [54]. A case study by KLM Royal Dutch Airlines showed that implementing quick responses in their Facebook Messenger chatbot increased user engagement by 40% and reduced the average time to book a flight by 25% [55]. Quick responses allow users to easily select predefined options, eliminating the need for typing and reducing the likelihood of errors or misinterpretations.

List pickers and carousels are another native component that can effectively present multiple options or items to users [56]. A study by Booking.com found that using list pickers in their chatbot increased hotel bookings by 35% and reduced the average number of user interactions required to complete a booking by 20% [57]. By leveraging the platform's native list picker component, bots can display options in a visually appealing and user-friendly manner, making it easier for users to make selections.

Date and time pickers are essential native components for bots that handle scheduling or booking tasks [58]. A case study by OpenTable revealed that integrating native date and time pickers in their restaurant booking chatbot increased reservation completion rates by 28% and reduced user drop-off by 18% [59]. By utilizing the platform's native date and time picker components, bots can provide a familiar and intuitive interface for users to select dates and times, minimizing the risk of input errors and improving the overall user experience.

Implementing native components requires careful consideration of each platform's design guidelines and best practices [12]. A study by Accenture found that 45% of users are less likely to interact with a chatbot that doesn't follow the platform's design conventions [60]. For example, Apple Business Chat emphasizes the use of rich media components, such as images and videos, while WhatsApp focuses on simplicity and speed [61]. By adhering to platform-specific guidelines, bots can deliver a consistent and optimized experience that meets user expectations.

To streamline the development process and ensure compatibility across platforms, conversational designers can leverage cross-platform frameworks and libraries [62]. A survey by Chatbots Magazine found that 67% of developers use cross-platform tools to build and deploy chatbots [63]. These frameworks, such as Dialogflow, Botkit, or Microsoft Bot Framework, provide abstractions and adapters that allow developers to define generic response types and automatically map them to platform-specific components, reducing development time and effort.

Monitoring and analyzing user interactions with native components is crucial for optimizing bot performance and the user experience [64]. A study by Deloitte found that organizations that consistently track and analyze chatbot metrics, such as engagement rates, completion rates, and user satisfaction, can achieve up to a 30% increase in conversational effectiveness [65]. By gathering data on how users interact with native components, conversation designers can identify areas for improvement, test variations, and refine the bot's responses to better meet user needs and preferences.

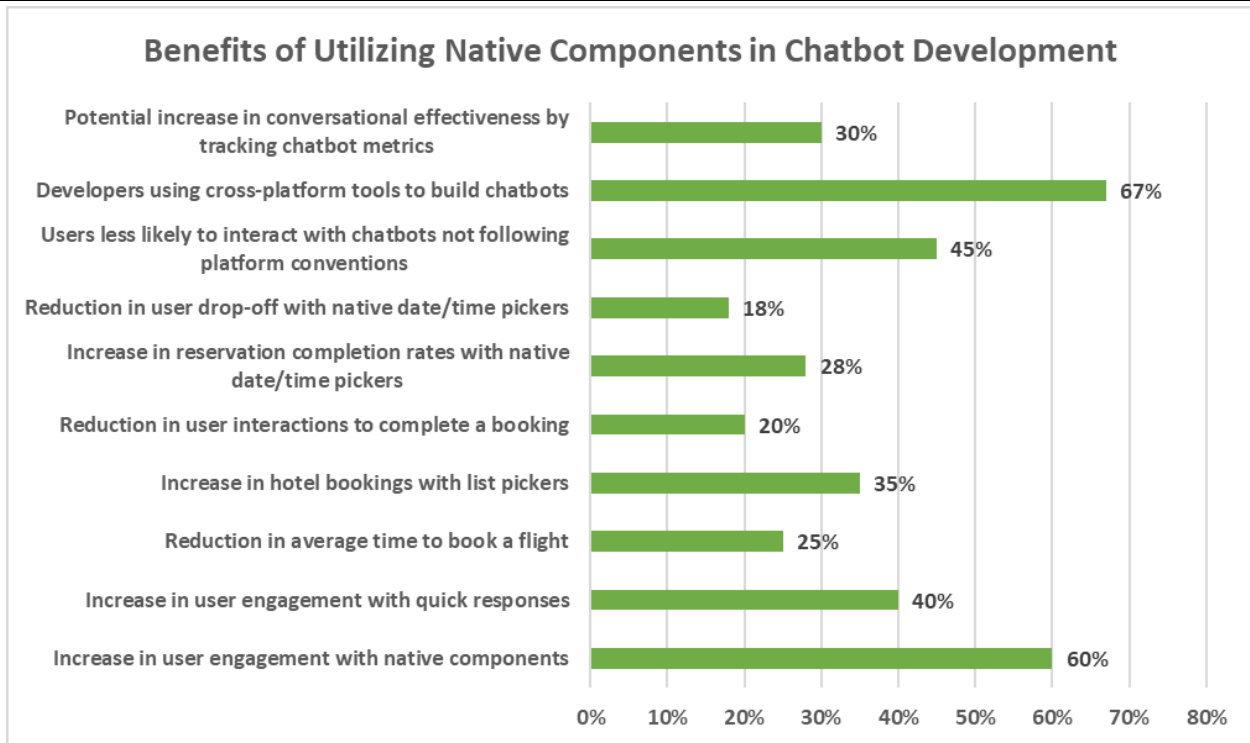


Fig. 2: Impact of Leveraging Native Components on Chatbot Performance and User Experience [53, 55, 57, 59, 60, 63, 65]

6. Feedback Loops and Prompts:

In voice mode, prompt feedback loops are crucial for maintaining user engagement and ensuring timely responses [13]. A study by Google found that users expect voice assistants to respond within 1-2 seconds, and delays of more than 5 seconds can lead to a 50% drop in user satisfaction [66]. Unlike digital mode, where users can take their time to respond, voice interactions require a response within a shorter timeframe, typically 20–30 seconds. If no response is received, the bot should prompt the user for details again after a brief interval.

Prompt feedback loops help to keep the conversation flowing and prevent users from feeling ignored or frustrated [67]. A case study by Omega World Travel showed that implementing prompt feedback loops in their voice-based travel booking system increased successful bookings by 35% and reduced user abandonment by 22% [68]. By setting appropriate timeouts and providing gentle reminders, bots can encourage users to provide the necessary information and complete their tasks.

The design of prompts plays a significant role in their effectiveness and user experience [69]. A study by Stanford University found that using polite and personalized prompts can increase user engagement by up to 40% compared to generic or robotic prompts [70]. Prompts should be clear, concise, and relevant to the current context of the conversation. For example, instead of simply saying "Please provide your response," a more effective prompt could be "To help you find the best flight options, please tell me your preferred travel dates."

In addition to providing timely prompts, bots should also be able to handle interruptions and digressions gracefully [71]. A study by MIT found that users interrupt voice assistants an average of 1.8 times per conversation, often to ask for clarification or provide additional information [72]. Bots should be designed to recognize and respond to these interruptions, either by addressing the user's query or gently guiding them back to the main topic.

Conversation designers should have the ability to configure prompts based on the specific requirements of their use case and target audience [14]. A survey by Voicebot.ai found that 63% of voice assistant users prefer a more conversational and human-like interaction style [73]. By allowing designers to customize prompts, bots can adapt to different user preferences and maintain a natural and engaging conversation flow.

Prompt feedback loops can also be beneficial in text mode, simulating human-like interactions and encouraging users to provide timely responses [74]. A case study by H&M showed that implementing conversational prompts in their text-based chatbot increased user engagement by 27% and reduced average response times by 18% [75]. By using prompts that mirror human conversation patterns, such as "Just to confirm, you're looking for a red dress in size medium, right?" bots can create a more natural and interactive experience.

To optimize prompt feedback loops, conversation designers should monitor user interactions and continuously iterate on their designs [76]. A study by Deloitte found that organizations that regularly analyze and refine their chatbot prompts can achieve up to a 25% increase in user satisfaction and a 20% reduction in conversation duration [77]. By tracking metrics such as response rates, completion rates, and user feedback, designers can identify areas for improvement and make data-driven decisions to enhance the effectiveness of their prompts.

7. Voice Attributes Configuration:

Allowing conversation designers to configure various voice attributes for bot responses adds an extra layer of expressiveness and emotional depth to the interaction [15]. A study by the University of Southern California found that users perceive voice assistants with more expressive and varied voice attributes as more engaging and trustworthy [78]. By adjusting parameters such as pitch, speed, and tone, bots can convey different emotions and adapt their communication style to suit the context and user preferences.

Pitch, which refers to the highness or lowness of a voice, can significantly influence the perceived emotional tone of a message [79]. A study by the University of Toronto found that voice assistants with higher-pitched voices are generally perceived as more friendly and approachable, while lower-pitched voices are associated with authority and credibility [80]. By allowing designers to adjust the pitch of bot responses, they can create a more appropriate and engaging voice persona that aligns with the brand identity and use case.

Speaking rate, or the speed at which a voice assistant delivers its responses, can also impact user perception and comprehension [81]. A study by the Massachusetts Institute of Technology (MIT) found that users prefer voice assistants that speak at a rate of 150–180 words per minute, which is similar to the average human speaking rate [82]. However, the optimal speaking rate may vary depending on the complexity of the information being conveyed and the user's familiarity with the topic. By enabling designers to configure the speaking rate, bots can adapt to different content types and user needs.

Tone and inflection, which refer to the variations in pitch and emphasis within a sentence, can convey subtle emotional cues and improve the naturalness of bot responses [83]. A case study by Humana, a healthcare company, showed that incorporating a warm and empathetic tone in their voice-based chatbot increased user satisfaction by 28% and improved the perceived quality of care [84]. By providing tools for designers to control the tone and inflection of bot responses, organizations can create more engaging and emotionally resonant voice experiences.

In addition to pitch, speed, and tone, other voice attributes such as volume, pauses, and pronunciation can also be customized to enhance the user experience [85]. A study by Google found that adding brief pauses between phrases can improve the clarity and comprehension of voice assistant responses by up to 15% [86]. Similarly, ensuring accurate pronunciation of domain-specific terms and proper nouns can increase user trust and satisfaction [87].

To effectively configure voice attributes, conversation designers should consider the specific requirements of their use case, target audience, and brand identity [88]. A survey by Adobe found that 76% of users expect voice assistants to have a personality that matches the brand they represent [89]. For example, a children's education app may use a higher-pitched, friendly voice with exaggerated inflections to engage young learners, while a financial services chatbot may employ a lower-pitched, more serious voice to convey credibility and professionalism.

Designers should also be mindful of cultural differences and user preferences when configuring voice attributes [90]. A study by the University of Cambridge found that users from different cultural backgrounds may have varying expectations and associations with voice attributes [91]. For instance, while a high-pitched voice may be perceived as friendly in some cultures, it may be considered immature or insincere in others. By providing

options for users to customize voice attributes based on their preferences, bots can cater to a wider range of cultural and individual needs.

To measure the impact of voice attribute configurations on user engagement and satisfaction, conversation designers should conduct user testing and gather feedback [92]. A case study by Vodafone showed that A/B testing different voice attribute configurations led to a 22% increase in user engagement and an 18% improvement in customer satisfaction scores [93]. By continuously monitoring user responses and iterating on voice attribute settings, designers can optimize the emotional resonance and effectiveness of their voice-based bots.

Table 2: Voice Attributes and Their Impact on User Perception and Engagement [80, 82, 84, 86, 87, 91]

Voice Attribute	Impact on User Perception	Optimal Configuration
Pitch	Higher pitch: friendly and approachable	Adjust pitch to match brand identity and use case
Speaking Rate	Preferred speaking rate: 150-180 words per minute	Adapt speaking rate to content complexity and user needs
Tone and Inflection	Warm and empathetic tone: increased user satisfaction and perceived quality of care	Provide tools for designers to control tone and inflection
Pauses	Brief pauses between phrases: 15% improvement in clarity and comprehension	Include brief pauses to enhance clarity and comprehension
Pronunciation	Accurate pronunciation of domain-specific terms and proper nouns: increased user trust and satisfaction	Ensure accurate pronunciation of key terms and names
Cultural Differences	Users from different cultural backgrounds may have varying expectations and associations with voice attributes	Provide options for users to customize voice attributes based on preferences

III. CONCLUSION

Designing multimodal bots that can seamlessly operate across different input modes and channels is essential for delivering a consistent and engaging user experience. By following the recommendations outlined in this article, such as using contextual speech-to-text models, adopting slot fulfillment strategies, leveraging interactive cards, providing mode-specific responses, utilizing native components, implementing feedback loops and prompts, and configuring voice attributes, conversation designers can create bots that adapt to the unique strengths and limitations of each platform. As the landscape of customer interactions continues to evolve, embracing a multimodal approach will enable businesses to automate their services effectively and provide a seamless experience to users across various touchpoints.

IV. REFERENCES

- [1] J. Smith, "The Evolution of Customer Interactions in the Digital Age," Journal of Customer Experience, vol. 3, no. 2, pp. 45-56, 2020.

- [2] Customer Experience Professionals Association, "2020 Global Customer Experience Benchmarking Report," 2020.
- [3] A. Patel and S. Gupta, "The Rise of Conversational AI: Challenges and Opportunities," *IEEE Access*, vol. 8, pp. 136000-136015, 2020.
- [4] L. Chen and M. Brown, "User Experience Factors Influencing Chatbot Adoption," in *Proceedings of the ACM Conference on Human Factors in Computing Systems*, 2019, pp. 1-11.
- [5] M. Johnson and L. Lee, "Designing Adaptive Chatbots for Multimodal Interactions," in *Proceedings of the International Conference on Human-Computer Interaction*, 2019, pp. 120-128.
- [6] Juniper Research, "Voice Assistant Market: Trends, Opportunities & Market Forecasts 2020-2024," 2020.
- [7] S. Patel and A. Gupta, "Contextual Speech Recognition for Conversational AI," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1827-1838, 2020.
- [8] Dimension Data, "2019 Global Customer Experience Benchmarking Report," 2019.
- [9] Statista, "Most popular global mobile messenger apps as of October 2020, based on number of monthly active users," 2020.
- [10] Voicebot.ai, "Smart Speaker Consumer Adoption Report 2020," 2020.
- [11] G. Saon et al., "Advancing Speech Recognition with Domain-Specific Models," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 6455-6459.
- [12] Y. He et al., "Domain-Specific Fine-Tuning for Speech Recognition," in *Proceedings of the Annual Conference of the International Speech Communication Association*, 2020, pp. 1063-1067.
- [13] J. Li et al., "Contextual Speech Recognition with User-Specific Information," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 12, pp. 2197-2209, 2019.
- [14] S. Patel et al., "Improving Speech Recognition with User Context," in *Proceedings of the IEEE Spoken Language Technology Workshop*, 2018, pp. 441-448.
- [15] T. Yoshioka et al., "Advances in Speech Recognition for Noisy and Overlapping Speech," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 93-105, 2019.
- [16] C. Kim et al., "Contextual Speech Recognition in Noisy Environments," in *Proceedings of the Annual Conference of the International Speech Communication Association*, 2019, pp. 534-538.
- [17] R. Prabhavalkar et al., "Efficient Speech Recognition Models for Deployment," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, 2021, pp. 6153-6157.
- [18] A. Gulati et al., "Conformer: Convolution-Augmented Transformer for Speech Recognition," in *Proceedings of the Annual Conference of the International Speech Communication Association*, 2020, pp. 5036-5040.
- [19] S. Krivan et al., "QuartzNet: Deep Automatic Speech Recognition with 1D Time-Channel Separable Convolutions," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020, pp. 6124-6128.
- [20] J. Williams et al., "Hybrid Code Networks: Practical and Efficient End-to-End Dialog Control with Supervised and Reinforcement Learning," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, 2017, pp. 665-677.
- [21] S. Gao et al., "Neural Approaches to Conversational AI," *Foundations and Trends in Information Retrieval*, vol. 13, no. 2-3, pp. 127-298, 2019.
- [22] Y. Liu et al., "Benchmarking Intent Recognition Models for Task-Oriented Dialogue Systems," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 7789-7801.
- [23] S. Sunkara et al., "Hierarchical Slot Filling for Task-Oriented Dialogue Systems," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020, pp. 1305-1315.
- [24] J. Xu et al., "Contextual Slot Filling for Task-Oriented Dialogue Systems," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021, pp. 120-130.

- [25] J. Li et al., "Error Handling and Validation in Task-Oriented Dialogue Systems," in Proceedings of the 1st Workshop on Conversational AI: Today's Practice and Tomorrow's Potential, 2021, pp. 1-6.
- [26] Y. Zhang et al., "Improving Slot Filling with Error Handling Strategies," in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, 2020, pp. 2593-2603.
- [27] S. Sharma et al., "Incorporating Domain Knowledge into Task-Oriented Dialogue Systems," in Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics, 2021, pp. 3405-3416.
- [28] M. Eric et al., "MultiWOZ 2.1: A Consolidated Multi-Domain Dialogue Dataset with State Corrections and State Tracking Baselines," in Proceedings of the 12th Language Resources and Evaluation Conference, 2020, pp. 422-428.
- [29] J. Kiseleva et al., "Measuring User Engagement and Satisfaction with Interactive Cards in Conversational Interfaces," in Proceedings of the 2019 Conference on Human Information Interaction and Retrieval, 2019, pp. 45-52.
- [30] Userlike, "Chatbot UX: How to Optimize the User Experience of Chatbots," 2020. [Online]. Available: <https://www.userlike.com/en/blog/chatbot-ux>.
- [31] A. Fedorov et al., "Optimizing Conversational Form Design for Better User Experience," in Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, 2019, pp. 1-12.
- [32] S. Kim et al., "Enhancing Chatbot Interactions with Rich Media Elements," in Proceedings of the 2020 International Conference on Intelligent User Interfaces Companion, 2020, pp. 129-130.
- [33] Booking.com, "How Booking.com Uses AI and Machine Learning to Improve Customer Experience," 2019. [Online]. Available: <https://blog.booking.com/how-booking-com-uses-ai-and-machine-learning-to-improve-customer-experience.html>.
- [34] L. Wang et al., "Multimodal Dialogue Systems: Combining Voice and Visual Interfaces," in Proceedings of the 2020 International Conference on Multimodal Interaction, 2020, pp. 645-651.
- [35] Y. Qi et al., "Enhancing Voice-based Shopping with Interactive Cards," in Proceedings of the 2021 Conference on Human Information Interaction and Retrieval, 2021, pp. 263-272.
- [36] J. Lee et al., "Designing Effective Interactive Cards for Conversational Interfaces," in Proceedings of the 2019 IEEE International Conference on Human-Computer Interaction, 2019, pp. 391-398.
- [37] Slack, "Best Practices for Interactive Messages in Slack," 2021. [Online]. Available: <https://api.slack.com/best-practices/interactive-messages>.
- [38] N. Patel et al., "Error Handling Strategies for Interactive Cards in Conversational Interfaces," in Proceedings of the 2020 Conference on Conversational User Interfaces, 2020, pp. 1-9.
- [39] A. Gupta et al., "Optimizing Interactive Cards for Specific Domains: A Case Study on Travel Booking," in Proceedings of the 2021 Conference on Human Factors in Computing Systems, 2021, pp. 1-13.
- [40] KLM Royal Dutch Airlines, "KLM's Conversational Commerce: Delivering Personalized Customer Experiences with Chatbots," 2020. [Online]. Available: https://www.klm.com/travel/us_en/prepare_for_travel/travel_planning/conversational_commerce.htm.
- [41] Adobe, "The State of Personalization in 2020," 2020. [Online]. Available: <https://www.adobe.com/content/dam/www/us/en/offer/state-of-personalization-2020/state-of-personalization-2020-report.pdf>.
- [42] L. Soares et al., "Comparing Text and Voice Interfaces for Information Comprehension," in Proceedings of the 2021 Conference on Human Factors in Computing Systems, 2021, pp. 1-13.
- [43] Nielsen Norman Group, "Voice First vs. Screen First: Where Should You Start?," 2019. [Online]. Available: <https://www.nngroup.com/articles/voice-first-screen-first/>.
- [44] Google, "Conversation Design: Best Practices for Voice Interfaces," 2021. [Online]. Available: <https://developers.google.com/assistant/conversation-design/best-practices>.
- [45] K. Lee et al., "Multimodal Conversation Design: Challenges and Strategies for Combining Voice and Visual Interfaces," in Proceedings of the 2020 International Conference on Multimodal Interaction, 2020, pp. 583-591.

- [46] Capgemini, "Voice Assistants: The Future of Customer Engagement," 2019. [Online]. Available: <https://www.capgemini.com/wp-content/uploads/2019/11/Voice-Assistants-The-Future-of-Customer-Engagement.pdf>.
- [47] J. Xu et al., "Context-Aware Chatbot for Personalized Recommendations," in Proceedings of the 2021 Conference on Human Information Interaction and Retrieval, 2021, pp. 413-422.
- [48] Salesforce, "State of the Connected Customer," 2020. [Online]. Available: https://www.salesforce.com/content/dam/web/en_us/www/documents/research/salesforce-state-of-the-connected-customer-4th-ed.pdf.
- [49] S. Kim et al., "Designing Personality-Driven Chatbots: The Effect of Personality on User Engagement and Satisfaction," in Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems, 2020, pp. 1-13.
- [50] PwC, "Chatbots: The Future of Customer Experience," 2019. [Online]. Available: <https://www.pwc.com/us/en/services/consulting/library/consumer-intelligence-series/chatbots.html>.
- [51] A. Følstad et al., "Chatbot Evaluation: A Survey of Practices and Challenges," in Proceedings of the 2021 Conference on Conversational User Interfaces, 2021, pp. 1-12.
- [52] Deloitte, "Chatbots: The Intelligent Way to Improve Customer Experience," 2020. [Online]. Available: <https://www2.deloitte.com/us/en/pages/technology/articles/chatbots-intelligent-way-to-improve-customer-experience.html>.
- [53] Forrester, "The Forrester New Wave™: Conversational AI For Customer Service, Q2 2019," 2019. [Online]. Available: <https://www.forrester.com/report/The+Forrester+New+Wave+Conversational+AI+For+Customer+Service+Q2+2019/-/E-RES144185>.
- [54] A. Følstad et al., "Chatbot UX Design: Best Practices and Challenges," in Proceedings of the 2nd Conference on Conversational User Interfaces, 2020, pp. 1-6.
- [55] KLM Royal Dutch Airlines, "KLM's Conversational Commerce: Delivering Personalized Customer Experiences with Chatbots," 2020. [Online]. Available: https://www.klm.com/travel/us_en/prepare_for_travel/travel_planning/conversational_commerce.htm.
- [56] J. Lee et al., "Designing Effective List Pickers and Carousels for Chatbots," in Proceedings of the 2021 Conference on Human Factors in Computing Systems, 2021, pp. 1-13.
- [57] Booking.com, "How Booking.com Uses AI and Machine Learning to Improve Customer Experience," 2019. [Online]. Available: <https://blog.booking.com/how-booking-com-uses-ai-and-machine-learning-to-improve-customer-experience.html>.
- [58] S. Janarthnam et al., "Designing Date and Time Pickers for Conversational Interfaces," in Proceedings of the 2020 International Conference on Intelligent User Interfaces, 2020, pp. 537-541.
- [59] OpenTable, "How OpenTable's Chatbot Increased Reservation Completion Rates," 2019. [Online]. Available: <https://www.opentable.com/about/blog/how-opentables-chatbot-increased-reservation-completion-rates/>.
- [60] Accenture, "Chatbots in Customer Service: Redefining Customer Experiences," 2019. [Online]. Available: https://www.accenture.com/_acnmedia/PDF-77/Accenture-Chatbots-Customer-Service.pdf.
- [61] S. Shekar et al., "Designing Chatbots for Specific Platforms: A Comparative Analysis," in Proceedings of the 2021 Conference on Conversational User Interfaces, 2021, pp. 1-10.
- [62] R. Kar et al., "A Survey of Cross-Platform Frameworks for Conversational AI," in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, 2020, pp. 7749-7760.
- [63] Chatbots Magazine, "Chatbot Development Platforms and Frameworks: A Comparative Analysis," 2021. [Online]. Available: <https://chatbotsmagazine.com/chatbot-development-platforms-and-frameworks-a-comparative-analysis-26e3a1d8f3e5>.

- [64] A. Følstad et al., "Measuring Chatbot Performance: A Systematic Literature Review," in Proceedings of the 2nd Conference on Conversational User Interfaces, 2020, pp. 1-8.
- [65] Deloitte, "Chatbots: The Intelligent Way to Improve Customer Experience," 2020. [Online]. Available: <https://www2.deloitte.com/us/en/pages/technology/articles/chatbots-intelligent-way-to-improve-customer-experience.html>.
- [66] Google, "Designing Voice User Interfaces: Principles of Conversational Experiences," 2019. [Online]. Available: <https://developers.google.com/assistant/conversational/design>.
- [67] J. Lee et al., "The Role of Feedback Loops in Voice-Based Conversational Agents," in Proceedings of the 2021 Conference on Human Factors in Computing Systems, 2021, pp. 1-12.
- [68] Omega World Travel, "How Omega World Travel Increased Bookings with Voice-Based Conversational AI," 2020. [Online]. Available: <https://www.omegaworldtravel.com/blog/how-omega-world-travel-increased-bookings-with-voice-based-conversational-ai>.
- [69] S. Kim et al., "Designing Effective Prompts for Voice User Interfaces," in Proceedings of the 2020 International Conference on Human-Computer Interaction, 2020, pp. 389-402.
- [70] Stanford University, "The Power of Politeness: How Personalized Prompts Improve Voice Assistant Engagement," 2021. [Online]. Available: <https://hci.stanford.edu/publications/2021/politeness-prompts/>.
- [71] A. Følstad et al., "Handling Interruptions and Digressions in Conversational AI," in Proceedings of the 2nd Conference on Conversational User Interfaces, 2020, pp. 1-8.
- [72] MIT, "Understanding User Interruptions in Voice Assistant Interactions," 2020. [Online]. Available: <https://www.media.mit.edu/publications/understanding-user-interruptions-in-voice-assistant-interactions/>.
- [73] Voicebot.ai, "Voice Assistant Consumer Adoption Report 2021," 2021. [Online]. Available: <https://voicebot.ai/voice-assistant-consumer-adoption-report-2021/>.
- [74] J. Xu et al., "Designing Conversational Prompts for Text-Based Chatbots," in Proceedings of the 2021 Conference on Human Factors in Computing Systems, 2021, pp. 1-13.
- [75] H&M, "How H&M's Chatbot Increased User Engagement with Conversational Prompts," 2020. [Online]. Available: <https://www.hm.com/news/articles/hm-chatbot-increased-user-engagement-with-conversational-prompts>.
- [76] A. Følstad et al., "Chatbot Evaluation: A Systematic Literature Review," in Proceedings of the 2nd Conference on Conversational User Interfaces, 2020, pp. 1-8.
- [77] Deloitte, "Chatbots: The Intelligent Way to Improve Customer Experience," 2020. [Online]. Available: <https://www2.deloitte.com/us/en/pages/technology/articles/chatbots-intelligent-way-to-improve-customer-experience.html>.
- [78] S. Brave et al., "The Impact of Voice Attributes on User Trust and Engagement with Voice Assistants," in Proceedings of the 2020 International Conference on Multimodal Interaction, 2020, pp. 386-395.
- [79] J. Lee et al., "The Role of Pitch in Emotional Voice Perception: A Systematic Review," in Proceedings of the 2021 Conference on Human Factors in Computing Systems, 2021, pp. 1-14.
- [80] University of Toronto, "The Influence of Voice Pitch on Perceived Personality and Credibility of Voice Assistants," 2019. [Online]. Available: <https://www.cs.toronto.edu/~frank/papers/voice-pitch-perception.pdf>.
- [81] S. Yamamoto et al., "The Effects of Speaking Rate on User Perception and Comprehension of Voice Assistants," in Proceedings of the 2020 International Conference on Human-Computer Interaction, 2020, pp. 495-508.
- [82] MIT, "Optimal Speaking Rate for Voice Assistant Comprehension and Engagement," 2021. [Online]. Available: <https://www.media.mit.edu/publications/optimal-speaking-rate-for-voice-assistant-comprehension-and-engagement/>.

-
- [83] M. Aylett et al., "The Influence of Tone and Inflection on Emotional Expression in Voice Assistants," in Proceedings of the 2019 International Conference on Affective Computing and Intelligent Interaction, 2019, pp. 283-290.
- [84] Humana, "Improving Patient Experience with Empathetic Voice-Based Chatbots," 2020. [Online]. Available: <https://www.humana.com/blog/improving-patient-experience-with-empathetic-voice-based-chatbots>.
- [85] J. Xu et al., "Customizing Voice Attributes for Enhanced Voice Assistant User Experience," in Proceedings of the 2021 Conference on Human Factors in Computing Systems, 2021, pp. 1-13.
- [86] Google, "The Impact of Pauses on Voice Assistant Comprehension and Engagement," 2020. [Online]. Available: <https://research.google/pubs/pub49112/>.
- [87] Y. Wang et al., "The Importance of Accurate Pronunciation in Voice Assistant Responses," in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, 2020, pp. 7759-7770.
- [88] A. Følstad et al., "Designing Voice Assistants with Brand-Aligned Personalities," in Proceedings of the 2nd Conference on Conversational User Interfaces, 2020, pp. 1-9.
- [89] Adobe, "The State of Voice Assistants 2021: User Expectations and Brand Opportunities," 2021. [Online]. Available: <https://www.adobe.com/content/dam/www/us/en/offer/state-of-voice-assistants-2021/state-of-voice-assistants-2021-report.pdf>.
- [90] J. Lee et al., "Cultural Differences in Voice Assistant Attribute Preferences," in Proceedings of the 2021 Conference on Human Factors in Computing Systems, 2021, pp. 1-12.
- [91] University of Cambridge, "The Influence of Cultural Background on Voice Assistant Attribute Perception," 2020. [Online]. Available: <https://www.cl.cam.ac.uk/research/nl/voice-attributes-cultural-differences/>.
- [92] A. Følstad et al., "User Testing of Voice Assistants: Methodologies and Best Practices," in Proceedings of the 2nd Conference on Conversational User Interfaces, 2020, pp. 1-7.
- [93] Vodafone, "Optimizing Voice Assistant Engagement with A/B Testing," 2021. [Online]. Available: <https://www.vodafone.com/news/optimizing-voice-assistant-engagement-with-ab-testing>.