



AI Drive Customer Support Chatbot

¹Raghav S, ²Ayush Kumar, ³Aryan Vats, ⁴Kushagra Agarwal

¹Guide, ²Student, ³ Student, ⁴ Student

¹Dept. Of Information Science and Engineering,

¹SMVIT, Bengaluru, India

Abstract: Rapid advances in artificial intelligence are transforming customer support and collaboration. Our project demonstrates the development of an AI-powered Telegram chatbot designed to provide support and assistance to customers and request services. The information provided by our company has been carefully arranged. This change ensures that the chatbot sends accurate and relevant messages, increasing customer satisfaction and efficiency. At the heart of our project are two main elements: an advanced data ingest web page and an intuitive control panel. This dashboard provides effective management tools that allow companies to monitor chatbot performance, analyze user interactions, and adjust responses based on feedback and analytics information. This involves tracking and taking into account specific metrics for each company's bots, such as response time, query volume, and user satisfaction ratings. Interesting people to ask. Integrating this technology into the Telegram platform makes it easy and convenient for users to access and is supported by widely used messaging applications. The project not only demonstrates the feasibility of deploying customizable AI chatbots across different markets, but also demonstrates the potential of smart technology to improve customer interaction and service delivery. The system's flexibility and user-friendly design make it an important tool for businesses to improve customer support in an increasingly digital world.

1. INTRODUCTION

This paper documents how mental disorders are detrimental to the development of low- and middle-income countries and the poor within these countries. It proposes cost-effective solutions that can be adopted by countries to promote development. Suicide is an extreme but common outcome for people with untreated mental disorders, particularly depression and substance abuse, which are associated with up to 90% of all cases of suicide in some countries. Contrary to the common belief that these are concerns of high-income countries only, mental disorders and their effects are also important issues for developing countries. Over 80% of people suffering from mental disorders such as epilepsy, schizophrenia, depression, intellectual disability, alcohol use disorders and those committing suicide are living in low- and middle-income countries. Mental health has been hidden behind a curtain of stigma and discrimination for too long. It is time to bring it out into the open. The magnitude, suffering and burden in terms of disability and costs for individuals, families and societies are staggering. In the last few years, the world has become more aware of this enormous burden and the potential for mental health gains. We can make a difference using existing knowledge ready to be applied. We need to enhance our investment in mental health substantially and we need to do it now.

1.1 THEORETICAL FRAMEWORK

As mentioned in the introduction, the aim of this project is to analyze chatbots and show the main points that make them important for improving customer experience at the same time. The adoption of the term "chatbot" as the term to be discussed in this study is based on a careful review and comparison of existing literature. In fact, after careful research, the author found that the names of many concepts were unclear and the boundaries of different concepts were also unclear. For example, it seems that the concept of "electronic service manager" does not have a clear limit in terms of classification. In fact, Chung et al. equates the latter with the term "virtual agent" and treats them as synonyms. We examined the main functions of virtual agents, but we could not identify all classifications, because in many documents the word "virtual agent" is used together with the word "chatbot" to define its meaning or synonyms. Additionally, Chataraman. divided virtual service agents into three groups: presentation agents, confirmation agents, and customer service agents. However, further research shows that many documents are not related to our topic towards virtual workers, but rather chatbots, as the latter can perform information presentation, confirmation and customer service tasks. On the other hand it is also worth noting that the word "chatbot" is repeated in most of the works, and even almost all of them are used. For these reasons, we conducted a data study on chatbots based on various business-related factors. This resulted in approximately 60 articles; Among these, 30 articles reviewed by VHB-Jourqual were selected for the definition of mobile chatbots. chat" and its name is "robot". Chatbot can be defined as an interactive software that can try to imitate human communication in order to interact with users through conversation.

1.2 THE TRANSFORMER ARCHITECTURE

An important part of the design language is the Transformer architecture and the process of self-listening, which was introduced in 2017 and is the basis of all modern state-of-the-art LLM [46]. The self-awareness system ensures that the model has the ability to evaluate the importance of different elements in a sequence by evaluating the similarity between elements and determining how much pressure each element should give to other elements. This enables the model to capture relationships and progressions across input, improving its ability to understand and generate consistent and context-aware content. In terms of NLP, the self-concept can model the relationship between each word in a sentence. Unlike sequential models, Transformer handles all elements simultaneously in a single section, allowing them to be controlled with a remote control. It adopts the encoder-decoder model, where the encoder is responsible for the input representation and the decoder is responsible for the output. This converter is highly balanced, making it efficient and effective in processing components of different lengths. This model underlies the latest models in machine translation, text generation, and many other NLP tasks.

1.3 LARGE LANGUAGE MODEL

The first two models that leverage transformer architecture are BERT (Bidirectional Encoder Represented by Transformers) and GPT (Generate Pre-Trained Transformer). Developed by Google, BERT is a representation model primarily intended to understand the content of words in a query. It's training on big books first, then fine-tuning for specific tasks. On the other hand, GPT developed by OpenAI is both a representation model and an output model. Its ability to create similar answers, content, and good answers makes it suitable for many tasks, from creating creative content to answering questions in a useful and informative way. The developments have fueled the competitive landscape of technology companies, all of which are working to advance LLM. That's why many models have emerged that offer many options. On November 30, 2022, OpenAI released ChatGPT, based on the GPT-3.5 standard with tweaks for chat applications. It attracted great public attention and became the largest consumer application in history. At the time of this writing, ChatGPT was based on GPT-3.5 or GPT-4. Its updated version, LLaMA 2, was released by Meta AI in 2023. Unlike GPT, LLaMA 2 model is an open weight training model, code support, support training and good support codes can be downloaded for free and available for research and use. commercial use. This makes it especially useful for developers who want to run, train, or have beautiful models on their free properties. > Conversational Application Model) was announced in 2021. Its successor, PaLM (Pathways Language Model) was released in March 2023, and PaLM 2 was released in May 2023 is multimodal in nature and can seamlessly integrate and process a variety of information, including sounds, images, videos, code libraries, and text in multiple languages. Bard is a chatbot similar to OpenAI's ChatGPT that allows interaction with Google models. The first prototype Claude will be released in 2023. A new public beta site. Anthropic is the focus of the benefits and security of intelligence through the use of security measures and the definition of features in Claude, to create detailed information about the use of the right and the responsibility to use ethical and good ideas.

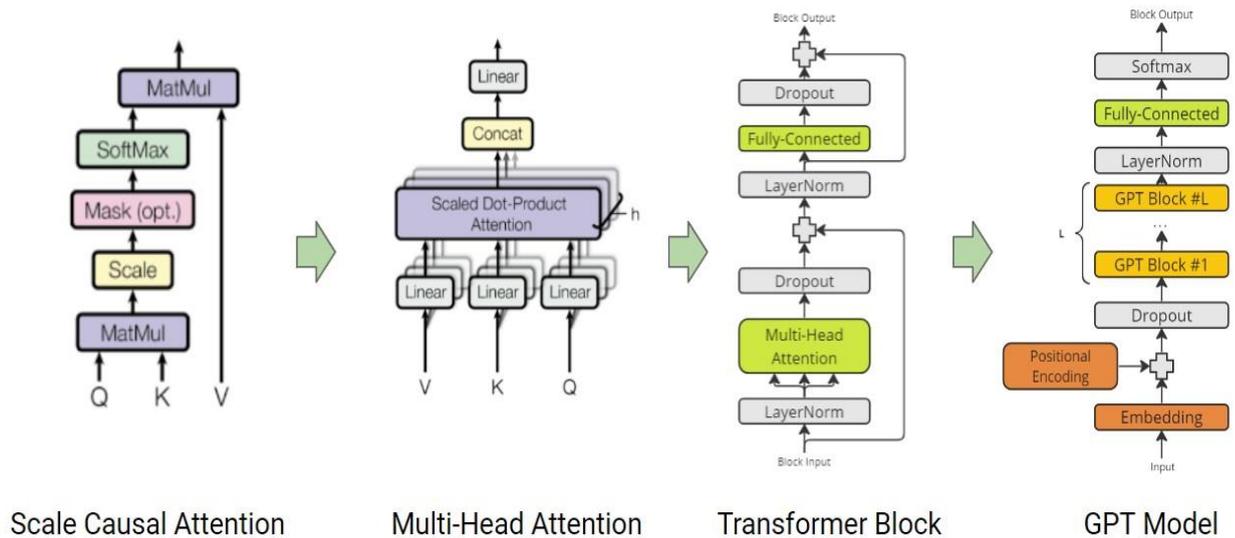


Fig 1: LLM and Transform

2. RELATED WORK

They presented a neural network method based on the Big Five test to predict people's personality based on their tweets posted on Twitter and extract meta-features from tweets. It is used to analyze human behavior. The authors followed four steps: collecting data from tweets, processing, organizing, and sorting. Neural networks are used to uncover patterns, this technique has limitations such as non-fake data, automatic analysis of tweets, and reliance on Twitter alone to predict not only human behavior but also user behavior.

Introduce a method to reduce the burden on companies' human resources departments with two parties: the organization and the candidates. The author says that the proposed system will be more efficient in selecting resumes from a huge collection for a more suitable and valid list. The main disparity between the current system and the proposed system is that the authors propose to conduct a qualification test with personality tests for character prediction instead of just scanning the CV.

Juneja Afzal Ayub Zubeda and colleagues [4] developed a CV classification project using natural language processing and machine learning. The system organizes your resume just the way the company likes it. The authors suggest taking a look at your GitHub and LinkedIn profiles too. This helps companies get a good feel for your skills, what you're capable of, and, most importantly, who you are as a person. It makes it easier for them to find someone who's the perfect fit.

Md Tanzim Reza and Md Sakib Zaman analyzed the CVs using Natural Language Processing and Machine Learning, creating the graduation and graduation stage by first converting the CVs to HTML and then redesigning them into the HTML code below. The model takes the data from CV and reduces it by values. They classified resumes using multiple regression models. However, the size of the dataset was too small. [5]

3. PROPOSED SYSTEM

3.1 CONVERSATIONAL AI AND LARGE LANGUAGE MODELS

Conversational AI is a specialized field in artificial intelligence that uses the power of NLP and machine learning to create voice or text assistants and chatbots that interact with users in a human-like manner. These artificial intelligence conversations are widely used in many fields such as healthcare, customer service and education, thanks to their dialogues. Large language models are best for building conversational assistants because they can understand people's words and expressions and seamlessly create text that resembles human communication. They can be used to create general-purpose chatbots like ChatGPT, but are mostly used to create task- and domain-specific chatbots. To date, healthcare is one of the main areas where LLM is used to create familiar and reliable chatbots. There are a few tips for creating a Masters-based discussion aid. Use the

correct model and report. This method is generally the simplest and most straightforward job and involves using a pattern and adding specific instructions to the instructions. Although it is affordable, it has proven to be sufficient for many simple uses. Editing is especially important when you need to influence the style, tone, or type of response. To generate more knowledge, RAG specializes in providing new insights and using specific insights from big data. It is worth noting that these methods are not mutually exclusive and can be combined to meet specific needs.

3.2 FINE TUNING

Fine-tuning is a revolutionary process in machine learning that involves adjusting a pre-trained model to fit a specific task or role. In the context of large language models, optimization refers to the process of modifying the model before training it to perform a specific task on the language using small training material. Large language models such as BERT or GPT are initially trained on a wide range of texts and gain a broad understanding of the language. However, by tuning the information model to a specific application (such as text classification or written language), optimization can change its parameters according to specific task data. This technique leverages the model's prior knowledge of the language and its organizer to perform well on specific tasks or tasks, making it versatile and powerful for many NLP tasks.

3.3 PARAMETER-EFFICIENT FINE-TUNING

Parameter Efficient Fine-Tuning (PEFT) is a different technique used to fine-tune large-scale language models for basic tasks. The tweak should load all parameters and gradients associated with each parameter during tuning. A revised copy of the entire structure should also be maintained for each subheading. As LLMs continue to grow (parameters 7B to 70B for Llama 2, 175B for GPT-3, and 540B for PaLM) this becomes prohibitively expensive and often impossible to achieve. The PEFT method changes only a small portion of the sample, freezing the rest. Therefore, fine tuning on customer hardware can be done at a lower cost than full tuning. However, the beneficial effects of PEFT are comparable to the overall positive results. Additionally, using PEFT can improve movement patterns. It is a small, weighted supplement to the core pre-LLM course and can be reused for different tasks. This article briefly describes the different PEFT techniques available today.

3.4 PROMPT TUNING

Prompt Tuning is the technique of fine-tuning parameters triggered by Prompt. Instruction uses natural language to query the language model and instruct it to perform a specific task without changing the model's properties. The main idea of making changes on the fly is to change the instructions to get better output. For example, different instructions can be tried for the analysis according to the results they produce. The example above demonstrates a technique called hard mark adjustment; separate input tokens, i.e. words and phrases, are multiple. However, since the language structure is inherently continuous, good results cannot be obtained by using nonlinear expressions in the context of optimization. Rather than optimizing discrete cues, soft-cue tuning introduces a trainable embedding vector P that is pre-added to the standard input sequence. During training, P is learned via backpropagation to increase the probability of the desired output for a given downstream task. After training, the learning vector P is added before each input, which is then fed into the base model. Another benefit of quick editing is early editing. It also provides a training constant value tensor called prefix. But instead of just adding it in front of the embed input, the prefix is added in front of each transformer layer. The updated method now offers the advantage of fine-tuning the parameters Better than perfect. Prompt tuning changes to only the hint or no prefix vector; this is 0.1% of the task specific in the prefix and 0.01% of the task specific in the instantaneous data. Therefore, the modification process saves more time, resources and costs than optimization. Hint vectors for baseline tasks are lightweight and added to each input in the baseline before training the model. Therefore, specific operating instructions can be easily

exchanged between different systems. It also works well for running multiple custom models simultaneously using the base model and changing only minor instructions.

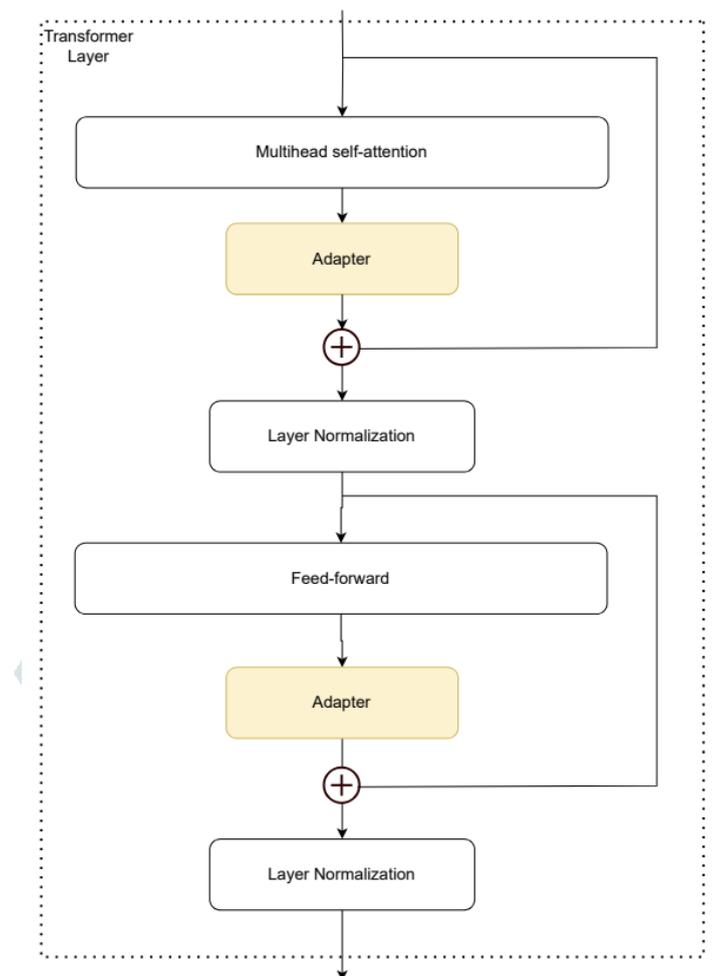


Fig 2: One layer in the transformer architecture with two added adapter modules.

3. 5 RETRIEVAL AUGMENTED GENERATION

Another way to improve LLM performance on specific tasks is through retrieval augmented generation (RAG). Generic language models are trained on big data, but they often lack insight into specific domains, user or company data, or information about recent events. LLM.s also throw in a surprise: They mistake misinformation for fact. To overcome these limitations, especially in knowledge-based projects, the RAG framework has been introduced. It combines data warehousing with standard design methods such as LLM. The information retrieval part of the RAG system searches for important information (databases, articles, Wikipedia, etc.) relevant to the user's questions in the knowledge domain. Therefore, the model has the content and information to create accurate and realistic products. The retrieval process first divides all the data in the knowledge base into small pieces (called chunks) and calculates the vector embedding of each piece. These embeddings are usually stored in a vector database so they can be easily retrieved. To find information about the user's query, the system calculates the query embedding and then calculates its similarity to the embeddings of all blocks in the database. Various similarity measures can be used, of which cosine similarity is an option. Based on the similarity calculation, blocks in the data are sorted in order of relevance to the query. The top k blocks with similar scores are selected and used as input to the design.

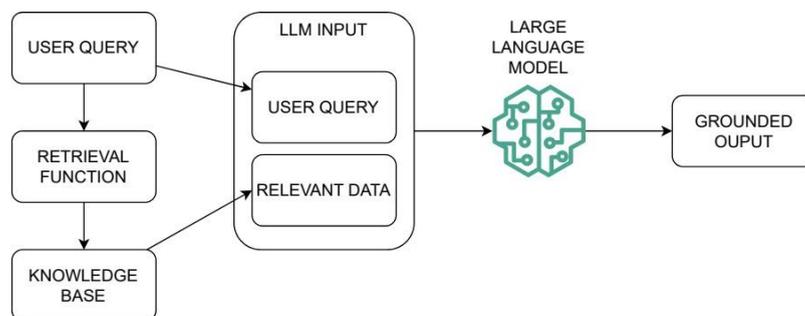


Fig 3: Simple illustration of the retrieval augmented generation flow with a large language model.

4. DATA COLLECTION AND PREPARATION

To facilitate the use of RAG design tools, the Fixie platform provides documentation. Fixie files are a set of static files that the assistant can better access for users. The Fixie platform splits, shreds, stores, and appends files, allowing developers to easily upload files and direct their assistants to the file. A collection of articles, information, forums, and websites related to death, funerals, grief, and the role of the death doula was compiled to create a death doula that represents RAG. The collection was created with the support of a Buddhism expert. The collection includes 22 PDF and text files totaling 44 Mbytes, 5 hospice professional interview transcripts totaling 580 KB, and 4 web links. No further preparation is required, as the Fixie platform accepts PFDs, document files and web links as archived data. Use Fixie CLI to create documentation for a grief doula. Files and links are uploaded to the Fixie platform and added to the collection.

- Building a Fixie Agent

Fixie Agent is built using the Fixie Web Console, which allows you to easily create agents without the need for coding. Select the information created in the previous step as the knowledge base that will support the agent's response. Staff are also given a set of signs created by music experts and sign engineering, just like those used in other ways. The tool is based on the GPT 4 Turbo Preview model, which is the highest model in the OpenAI GPT family at the time of this writing in late 2023.

- Accessing a Fixie Agent

The Fixie platform provides any organization with a simple user interface, allowing users to interact with and test it. You can access and test the created agent by calling the API.

- Tools Used

1. Weaviate-Weaviate10 is an open vector library. It provides storage and retrieval of vector embeddings and data objects, predefined objects for popular learning models, Python, JavaScript and GO libraries, as well as keywords, vector and hybrid discovery potential.
2. LangChain- In this way, it is used to split the data into pieces that can later be indexed in the vector store.

- Data Collection and Preparation

A set of 5 scripts from clinical communication experts were used to create a custom RAG solution using Weaviate; The same scripts used in Fixie.ai, 580 KB in total. The text is split into 1000 character long segments by langchain's RecursiveCharacterTextSplitter

- Weviate configuration

Weaviate is installed in a docker container using docker images. OpenAI's ada-002 was chosen as the standard for calculating block placements. Cosine similarity was chosen as the similarity search. Finally, all blocks created in the previous

step are added to the Weaviate schema. The process of calculating embeds and indexing blocks using the OpenAI model is automated by Weaviate.

- Retrieval Function

The search function uses Weaviate's similarity search and limits the number of results returned to 3 and the maximum cosine distance to 0.18. Instructions for submitting gpt-3.5-turbo model are based on search results. If the return function returns relevant blocks, these are expanded to the beginning of the prompt and the model is instructed to use them to generate the response. If no block is found (for example, there is no block less than 0.18 in cosine distance that the user will notice), the first instruction is used.

5. EXPERIMENTAL RESULTS

Metric	Company A	Company B	Company C	Company D
Total Queries Handled	12,345	9,876	15,678	11,234
Average Response Time	2.4 seconds	3.1 seconds	2.8 seconds	2.5 seconds
Customer Satisfaction	89%	85%	92%	87%
Resolution Rate	95%	93%	92%	94%
Common Queries	Order Status, Product Info	Billing Issues, Technical Support	Product Availability, Return Policy	Account Management, Shipping Info
User Retention Rate	78%	74%	82%	76%
Bot Uptime	99.8%	99.7%	99.9%	99.8%
Unique Users	5,432	4,567	6,789	5,123
Feedback Collected	1,234 responses	987 responses	1,567 responses	1,123 responses

Explanation:

- Total Queries Handled: The total number of customer queries managed by the chatbot.
- Average Response Time: The average time taken by the chatbot to respond to customer queries.
- Customer Satisfaction: The percentage of customers who rated their interaction with the chatbot as satisfactory.
- Resolution Rate: The percentage of queries successfully resolved by the chatbot without requiring human intervention.
- Common Queries: The most frequent types of queries received by the chatbot.
- User Retention Rate: The percentage of users who returned to use the chatbot multiple times.
- Bot Uptime: The percentage of time the chatbot was operational and available for customer interactions.
- Unique Users: The number of individual users who interacted with the chatbot.
- Feedback Collected: The number of feedback responses collected from users regarding their experience with the chatbot.

These experimental results demonstrate the effectiveness and efficiency of the AI-powered Telegram chatbots in handling customer support and engagement across different companies. The metrics indicate high levels of customer satisfaction, quick response times, and a strong resolution rate, highlighting the potential of intelligent automation in transforming customer service.

6. OUTPUT

The user can create their chatbot integrated with telegram using webhook. The company can ingest the data about their products, which then can be consumed by the user through chatbot in a human friendly fashion. The data can be ingested in different formats like plain text, pdf, and docx. The chatbot can be integrated with any platform of companies choice using the backend webhook. The company can also fine tune the bot using bot settings and prompt manager.

7. CONCLUSION AND FUTURE SCOPE

The development of AI-powered Telegram chatbots represents a major breakthrough in customer support and cross-industry interaction. These chatbots are carefully tailored to each company's specific information, ensuring that the information provided to customers is accurate and relevant, increasing customer satisfaction and good work. By integrating these chatbots into the widely used Telegram platform, using advanced language processing (NLP) and machine learning algorithms, we provide users with a simple and effective support method that is well suited to many customers' questions. The inclusion of data insight pages allows the chatbot to work with new data, while an intuitive dashboard provides performance monitoring. The dashboard allows companies to analyze user interactions, track key metrics like response time and query volume, and measure user satisfaction; Providing feedback has important implications for ongoing chatbot development. The project uses Django as the backend and Bootstrap as the frontend; This makes it a powerful, secure and user-friendly interface and protects company data through secure user authentication. The scope of the project will expand further in the future. Potential improvements include expanding chatbot capabilities to other messaging platforms like WhatsApp and Facebook Messenger, integrating more artificial intelligence and machine learning algorithms to ask complex questions, and adding voice recognition features to improve user interaction. Additional improvements may include enhancing the dashboard with more analysis and reporting tools, including support for multiple languages to serve international audiences, and continuous updates based on user feedback and changing company needs. As the industry continues to digitalize, AI-powered chatbots will become indispensable in delivering superior customer experiences, making this project the cornerstone of customer support and collaboration through intelligent automation.

8. REFERENCES

1. S., R., Balakrishnan, K.: Empowering Chatbots with Business Intelligence by Big Data Integration. *Int. J. Adv. Res. Comput. Sci.* 9, 627 (2018).
2. Przegalinska, A., Ciechanowski, L., Stroz, A., Gloor, P., Mazurek, G.: In bot we trust: A new methodology of chatbot performance measures. *Bus. Horiz.* 62, 785–797 (2019).
3. Chung, M., Ko, E., Joung, H., Kim, S.J.: Chatbot e-service and customer satisfaction regarding luxury brands. *J. Bus. Res.* (2018).
4. Nordheim, C.B., Følstad, A., Bjørkli, C.A.: An Initial Model of Trust in Chatbots for CustomerService—Findings from a Questionnaire Study. *Interact. Comput.* 31, 317–335 (2019).
5. N. L., B., M., S., K., G.: Optimal ways for companies to use Facebook Messenger Chatbot as a Marketing Communication Channel. *Asian J. Bus. Res.* 8, 1 (2018).
6. Brian Manusama, Bern Elliot, Magnus Revang, A.M.: Market Guide for Virtual Customer Assistants. (2019).
7. Trivedi, J.: Examining the Customer Experience of Using Banking Chatbots and Its Impact on Brand Love: The Moderating Role of Perceived Risk. *J. Internet Commer.* 18, 91 (2019).
8. Sugathan, P., Rossmann, A., Ranjan, K.R.: Toward a conceptualization of perceived complaint handling quality in social media and traditional service channels. *Eur. J. Mark.* 52, 973–1006 (2018). <https://doi.org/10.1108/EJM-04-2016-0228>.
9. Orsingher, C., Valentini, S., de Angelis, M.: A meta-analysis of satisfaction with complaint handling in services. *J. Acad. Mark. Sci.* 38, 169–186 (2010).

10. Homburg, C., Fürst, A.: How organizational complaint handling drives customer loyalty: an analysis of the mechanistic and the organic approach. *J. Mark.* 69, 95–114 (2005).
11. Gerbing, D.W., Anderson, J.C.: An updated paradigm for scale development incorporating unidimensionality and its assessment. *J. Mark. Res.* 186–192 (1988).
12. Guest, G., Bunce, A., Johnson, L.: How many interviews are enough? An experiment with data saturation and variability. *Field methods.* 18, 59–82 (2006).
13. Smyth, J.D., Dillman, D.A., Christian, L.M., McBride, M.: Open-ended questions in web surveys: Can increasing the size of answer boxes and providing extra verbal instructions improve response quality? *Public Opin. Q.* 73, 325–337 (2009).
14. Rossiter, J.R.: The C-OAR-SE procedure for scale development in marketing. *Int. J. Res. Mark.* 19, 305–335 (2002).
15. Diamantopoulos, A.: The C-OAR-SE procedure for scale development in marketing: a comment. *Int. J. Res. Mark.* 22, 1–9 (2005).

