



## Article

# Alquist 5.0: Dialogue Trees Meet Generative Models, a Novel Approach for Enhancing SocialBot Conversations

Ondrej Kobza <sup>\*,†</sup>, David Herel <sup>†</sup> , Jan Cuhel <sup>†</sup>, Tommaso Gargiani <sup>†</sup>, Petr Marek <sup>†</sup> and Jan Sedivy <sup>\*</sup>

Faculty of Electrical Engineering, Czech Technical University in Prague, Jugoslávských partyzánů 1580, 160 00 Praha, Czech Republic; david.herel@seznam.cz (D.H.); cuheljan@fel.cvut.cz (J.C.); tommaso.gargiani@gmail.com (T.G.)

\* Correspondence: kobzaond@fel.cvut.cz (O.K.); jan.sedivy@cvut.cz (J.S.); Tel.: +420-739-843-056 (O.K.)

<sup>†</sup> These authors contributed equally to this work.

**Abstract:** This article introduces Alquist 5.0, our SocialBot that was designed for the Alexa Prize SocialBot Grand Challenge 5. Building upon previous iterations, we present the integration of our novel neural response generator (NRG) Barista within a hybrid architecture that combines traditional predefined dialogues with advanced neural response generation. We provide a comprehensive analysis of the current state-of-the-art NRGs and large language models (LLMs), leveraging these insights to enhance Barista's capabilities. A key focus of our development was in ensuring the safety of our chatbot and implementing robust measures to prevent profanity and inappropriate content. Additionally, we incorporated a new search engine to improve information retrieval and response accuracy. Expanding the capabilities of our system, we designed Alquist 5.0 to accommodate multimodal devices, utilizing APL templates enriched with custom features to deliver an outstanding conversational experience complemented by an excellent user interface. This paper offers detailed insights into the development of Alquist 5.0, which effectively addresses evolving user demands while preserving its empathetic and knowledgeable conversational prowess across a wide range of topics.

**Keywords:** Alexa Prize; SocialBot; neural response generator; Barista; hybrid dialogues; safety; Vicuna; BlenderBot 3



**Citation:** Kobza, O.; Herel, D.; Cuhel, J.; Gargiani, T.; Marek, P.; Sedivy, J. Alquist 5.0: Dialogue Trees Meet Generative Models, a Novel Approach for Enhancing SocialBot Conversations. *Future Internet* **2024**, *16*, 344. <https://doi.org/10.3390/fi16090344>

Academic Editor: Gianluigi Ferrari

Received: 12 July 2024

Revised: 16 September 2024

Accepted: 18 September 2024

Published: 21 September 2024



**Copyright:** © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

This article delineates the development of Alquist 5.0, our SocialBot designed for the fifth iteration of the Alexa Prize SocialBot Grand Challenge 5 [1]. This university-based contest aspires to foster the creation of an Alexa skill that exhibits proficiency in coherent and engaging discourses about current affairs and common topics with humans for a sustained duration of 20 min, thereby aiming at achieving a median rating of at least 4.0 on a 1.0–5.0 scale. This is a challenging task requiring the integration of key dialogue components into our SocialBot. As such, we concentrated on developing an Alquist version that offers an empathetic, as well as knowledgeable, interaction on an array of subjects.

Alquist 5.0 is an evolution of its predecessors, and it was brought to life in earlier editions of the SocialBot Grand Challenge. Despite this, the latest version came into the scene at a time when the domain of conversational artificial intelligence witnessed a significant enhancement of its capabilities. This phenomenon is attributed to the enlargement of user expectations. The launch of OpenAI's ChatGPT (<https://openai.com/blog/chatgpt>, accessed on 30 November 2022) in November 2022 managed to spark widespread fascination in the arena of large language models (LLMs) by demonstrating its astounding potential to the general populace. Consequently, this has led to a wide gap that has to be covered by SocialBots vying for the Alexa Prize regarding meeting user expectations. This predicament is reflected in the decreased user-assigned ratings that have arisen relative to previous contests. Contrary to restricted-access, closed-source models such as ChatGPT

and Google's Bard (<https://bard.google.com/>, accessed on 2 October 2023), a myriad of open-source LLMs boasting performance analogous to their commercial counterparts have been released. This period of flourishing language modeling has facilitated us to pivot our focus onto constructing a novel neural response generator (NRG)—Barista (refer to Section 4)—for our SocialBot.

Incorporating Barista into Alquist has been carried out in a way such that the resulting system is flexible enough to handle various user responses well and to dive deep into various topics (i.e., each conversation can be different and there are just few repetitions across sessions), but it also has the ability to conduct qualitatively well-written dialogues (fine-grained and expanded dialogues based on the previous iteration, i.e., Alquist 4.0 [2]) for most common topics. Therefore, we pioneered several novel methods for integrating Barista into our SocialBot, leading to the enhancement of our SocialBot's generalization capabilities. A more detailed discussion on the topic is made available in Section 7.

Furthermore, we also designed our SocialBot for multimodal devices such as Alexas equipped with a touchscreen. We put to use the APL (<https://developer.amazon.com/en-US/docs/alexa/alexa-presentation-language/add-visuals-and-audio-to-your-skill.html>, accessed on 28 November 2023) templates, which were enriched with our custom feature implementations, with the ambition of matching an outstanding conversational experience with an equally excellent user interface. For more insights into this, refer to Section 8.

## 2. Related Work

The existing research dealing with open-domain chatbots primarily targets two avenues: the development of more advanced neural generation methods [3,4] (often achieved by integrating knowledge bases into generative models); and the established methodologies that typically rely on dialogue trees ([5]). An exemplar of the latter approach can be found in the architectural blueprint proposed by [6].

A major drawback of traditional dialogue tree systems is their limited scope. These systems rely heavily on predefined rules and scripts, which restrict their ability to handle the vast array of possible user inputs and conversational contexts. This rigidity often results in unnatural and repetitive interactions as the system cannot adapt to new or unexpected user queries. Additionally, the manual effort required to design and update these dialogue trees is considerable, making them less scalable for complex applications.

To address these limitations, several researchers have sought to enhance these predetermined rule systems by integrating statistical generative models [2,7].

Traditional generative models, like n-grams, have been replaced by more advanced neural network architectures as computational power has increased. Typically, these newer models use encoder–decoder networks [8] with LSTM [9] or other recurrent units. While these models have significantly improved language modeling and other NLP tasks, their design does not allow for high parallelization, thus limiting their size. A major breakthrough came with the introduction of the Transformer architecture [10], which forms the basis of most current state-of-the-art language models. This led to the development of new static models, such as BERT [11] and ALBERT [12], as well as advanced generative models, such as GPT [13] or T5 [14].

The prior state-of-the-art generative models, such as DialogueGPT [15] and T5 [14], had limited capacity to manage extensive, in-depth conversations. These models often struggled with maintaining context over long dialogues, leading to issues such as generating hallucinations—responses that are plausible-sounding but factually incorrect—and contradictions within the conversation. Furthermore, these models require substantial computational resources for training and deployment, which poses additional challenges for practical implementation. Hence, [2] decided to use traditional, manually created dialogues by professionals and to use their own version of DialogueGPT for out-of-domain inputs.

In November 2022, OpenAI (<https://openai.com/>, accessed on 2 December 2023) unveiled GPT3.5 and, subsequently, introduced GPT4. These two models represented a significant step forward for the natural language processing (NLP) field and effectively

initiated a new phase in the development of highly competent instruction-based large language models. The Llama model proposed by [16] and the Falcon LLM proposed by [17] can be employed in open-domain chatbots either as standalone systems (such as Chatbotarena (<https://stablediffusion.fr/vicuna>, accessed on 27 May 2023) or as part of an integrated pipeline. More recently, FastChat (<https://github.com/lm-sys/FastChat>, accessed on 30 March 2023) introduced the Vicuna model, a fine-tuned version of the Llama model optimized for an instruction-based dataset. In our study, we used the 7 billion-parameter Vicuna model, a 3 billion-parameter version of BlenderBot 3, and two other smaller language models that were integrated into our Barista NRG (DeBERTa-xsmall [11,18]) for classification, as well as FLAN-T5 [19] for query generation when external information was required. Regarding external information, we leveraged Fusion-in-Decoder [20]; however, other techniques such as Retrieval Augmented Generation [21] or a knowledge graph encoder [22] are also an option.

While dialogue safety issues can easily be controlled in dialogue tree systems by design, such challenges must be meticulously addressed when using large language models. The subject of chatbot dialogue safety has drawn considerable research attention in recent years, with studies focusing on identifying and mitigating the risks associated with unsafe language generation in conversational models [23–25]. For example, researchers have utilized the Wikipedia Toxic Comments dataset [26], which incorporates 100k pieces of human-labeled content, and a Twitter dataset [25], consisting of 240k tweets categorized as containing hate speech or offensive language by human annotators. Moreover, the literature has delved into other dialogue safety issues, including political prudence [27], stereotyping [28], and various de-biasing techniques for specific domains [24,29,30].

Although our research primarily concentrated on an NLP application (an open-domain chatbot more specifically), we also allocated substantial attention to our user interface. Ref. [31] demonstrated the significant impact that user interface design can have on the overall user experience. This inspired us to devise an outstanding user interface and unveil our novel digital persona—Alquistyna.

### 3. System Architecture

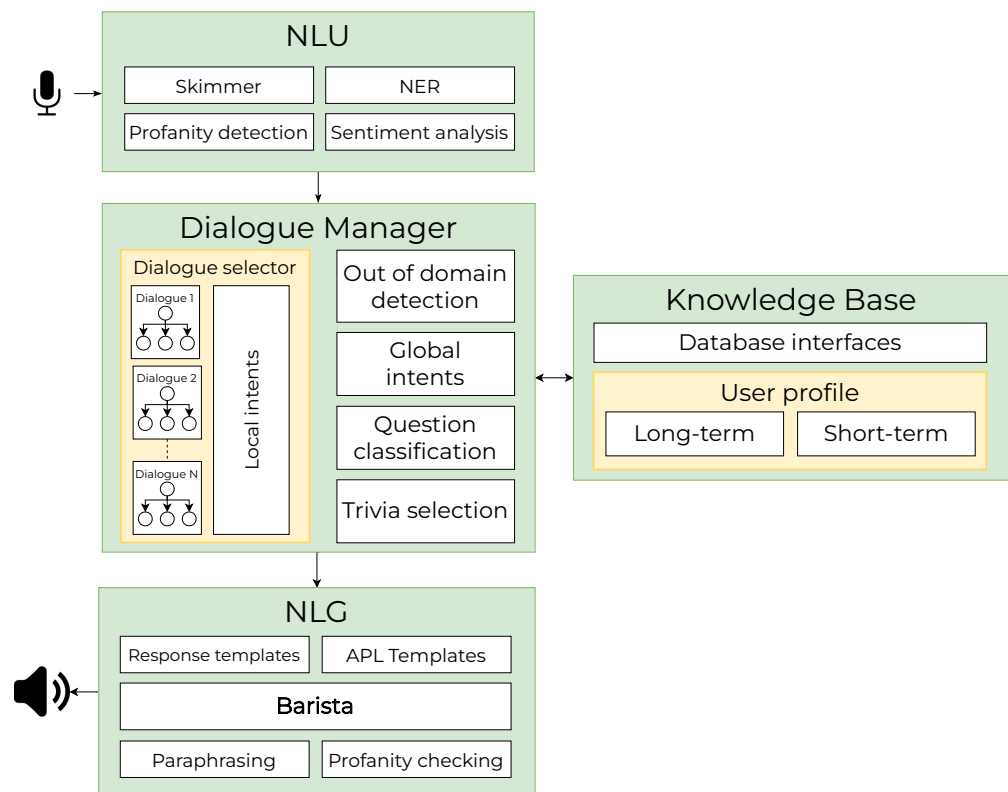
The system is based on the architecture proposed by [2] (see Figure 1). The main improvement we made for the architecture was its extension with the NRG Barista (see Section 4) together with its integration within traditional dialogue trees; thus, a hybrid architecture was built. Barista is used for out-of-domain conversations and in our hybrid dialogues, where several nodes in a dialogue tree are generated dynamically by our NRG. Barista is a neural response generator designed for open-domain conversations that access external knowledge sources to provide knowledgeable responses. We utilize Barista in various conversational contexts, as described in Section 7.

The invention of the hybrid architecture had several motivations:

1. In the previous iteration of the Amazon Alexa Prize SocialBot Grand Challenge (4) (<https://www.amazon.science/alex-prize/socialbot-grand-challenge/2020>, accessed on 14 August 2021), the best-performing SocialBots were generally based on manually prepared dialogues with some ability in handling out-of-domain cases by a neural response generator. No SocialBot that was based entirely (or almost entirely) on large language models has been provided with an award (although some SocialBots have had this architecture, e.g., Genuine by [32]).
2. Furthermore, some of the current chatbots solely based on LLMs, such as Meena [33] or BlenderBot ([4,34–36]), still have several qualitative problems [37,38]. We assume that carefully designed dialogues (with consultation help from linguists from PromethistAI (<https://promethist.ai/>, accessed on 1 January 2022)) for most common topics (hobbies, movies, food, books, music, etc.) will perform better (in the particular domains) as the dialogues should produce sensible and even interesting engaging outputs. This is not guaranteed for general LLMs.

3. Using only handcrafted dialogue trees poses huge limitations regarding generalization. Moreover, returning users might experience that conversations with such rule-based systems are (in longer periods) repetitive, which we assume would have a significant negative impact on the user experience. Incorporating our Barista NRG allows our SocialBot to converse about any topic. Furthermore, hybrid dialogues ensure within-dialogue variability.

Hence, the ratio behind the architecture choice is to combine the advantages of traditional rule-based chatbots with the advantages of LLM-based conversational systems.



**Figure 1.** The system architecture is built upon the work proposed by Konrád et al. [2], with the main emphasis being put on the Barista neural response generator.

From the dialogue design perspective, the main building block of the proposed SocialBot is the *main dialogue*. Within the main dialogue, we have a collection of dialogue groups composed of multiple dialogues. Each group is dedicated to a specific topic, such as sports, movies, or music. This design allows our SocialBot to engage in conversations about a wide range of subjects, providing a comprehensive and diverse user experience. Furthermore, each dialogue follows a tree graph structure. This means that we have the flexibility to specify a conversational flow, defining how the system should respond to different user utterances. By organizing the dialogue in this structured manner, we can ensure a coherent and intuitive conversation experience for the users.

In order to effectively understand user utterances, we followed the practice established by [2]. We classified every utterance into an in-domain or out-of-domain intent. In-domain intents are designed to handle user utterances that are anticipated within the context of the dialogue. Furthermore, we distinguished between two types of in-domain intents: *local*, which are valid only in a specific context and cover the expected user inputs that align with the predefined conversational flow; and *global*, which are valid across all dialogues and can be relevant in any part of the conversation. On the other hand, out-of-domain intents handle unexpected user utterances that fall outside the predefined scope of the dialogue. Out-of-domain intents play a crucial role when incorporating our new NRG Barista into the system. We also employed Skimmer (as proposed by [2]), which analyzes

user utterances and identifies relevant information. For example, if a user mentions that they have a dog, Skimmer can detect and extract this information without asking the user about it directly. This ability to gather user details implicitly helps in building the user profile, which is designed to store information about the user; this is useful for conversation (e.g., that they own a dog) and can be leveraged to provide personalized responses in subsequent interactions.

We navigate the user between the dialogue trees utilizing a component called Dialogue Selector, which chooses the most relevant topic for conversation. To achieve this, we applied the policy proposed by [2]. This selection is based on two factors. The first factor is the user profile, which includes details gathered through Skimmer about the user's preferences, hobbies, and other relevant information. This enables our SocialBot to suggest topics that align with the user's interests. The second factor is the topics identified during the previous interaction with the bot. These topics, detected using entity recognition, are then selected as possible candidates for conversation.

Most of our dialogues include questions for the user, and we have found that users sometimes respond by answering and then asking the same question back (e.g., "Bot: How are you?" "User: Fine, and you?"). Unfortunately, our dialogues do not have pre-written answers for such questions. To address this, we implemented a question classifier that detects these cases. When this happens, Alquist uses its neural response generator to answer the question, and then the next dialogue is selected.

Furthermore, we enrich the dialogues by inserting fun facts (i.e., interesting trivia information) between dialogues. Some dialogues also use external knowledge bases, specifically Amazon Evi (a QA engine that can be asked questions in a natural language, and it can be accessed via Cobot Toolkit [39]) for question answering and The Movie Database (<https://developer.themoviedb.org/docs>, accessed on 6 May 2023) (TMDB) for providing information about movies that are mentioned during the conversation.

#### 4. Barista Neural Response Generator

We introduce our innovative neural response generator (NRG) due to its versatility in different dialogue contexts and to improve handling out-of-domain conversations. Our approach is inspired by the BlenderBot 3 (BB3) chatbot [4]—we modified the original pipeline, reduced latency, and improved the overall performance. In its architecture, Barista uses two generative models:

1. A 3-billion-parameter model based on BlenderBot 3;
2. A 27-billion-parameter model based on Vicuna.

Figure 2 illustrates the architecture of Barista. The pipeline begins with a classification task to determine whether the NRG should query external knowledge. Concurrently, the system generates a search query for our search API (see Section 5). Suppose the output of the classifier is 1 ('yes'). In such a case, Barista retrieves information through our search API. In terms of search, the pipeline extracts specific information from a set of documents returned by the API. Subsequently, all external information, user input, and conversational context are fed into the Barista 3B BB3-based model, and the final utterance is then generated.

In parallel, Barista runs VicuChat NRG, our fine-tuned version of Vicuna 7B. VicuChat generates a response according to the user input, context, and possibly knowledge (together with the query), and it then generates the next utterance. Finally, a rule-based classifier decides which generated output to use.

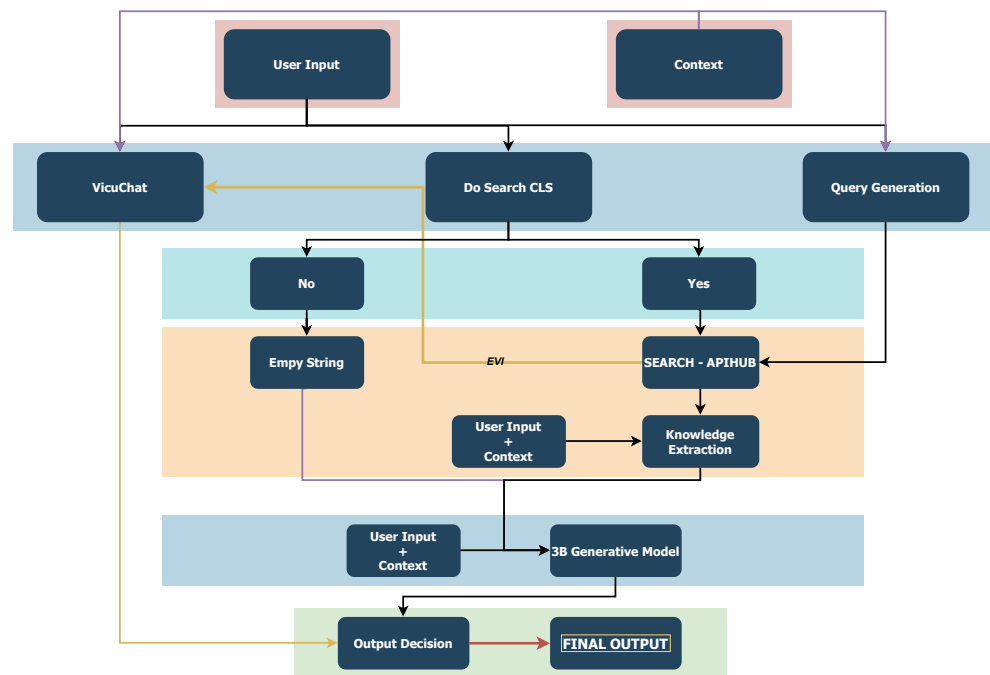


Figure 2. Visualization of Barista.

#### 4.1. Models

This section offers a concise review of the proposed models, placing a significant emphasis on the intricacies of the 3B and Vicuna-centric generative models. Additionally, it provides a comprehensive analysis of their implementation and potential applications. Regarding task-specific models, we took inspiration from the work published by [38].

##### 4.1.1. Conducting Search Classification and Query Generation

These modules inform the pipeline on whether to query external knowledge and, if yes, to generate the query. We trained our custom models, as proposed by [38] (DeBERTa-xsmall [18] for the classification task and FLAN-T5 [19] for query generation) on the same datasets that were used for BB3. The primary purpose was to reduce latency (our models are 17× faster in terms of classification, 4× faster in terms of query generation, and are achieving similar or higher performance [38]); however, usage of the BB3 model would have also been possible.

As [38] described (we used the same data for training), the training dataset for classification regarding search contains 490,000 samples, while the training dataset for query generation is composed of information relevant to conditional language modeling and contains ca. 40k samples. We trained the models on a single NVIDIA A100 GPU, and the training took only a few hours for both models.

##### 4.1.2. Knowledge Extraction

BB3’s knowledge extraction task involves extracting relevant information, typically short phrases or single sentences, from a set of documents using the FiD architecture ([20]). We used the same data as [4] (12 GB dataset for conditional language modeling) for this task and trained two DeBERTa base models. The first model was a classical encoder Transformer, which was proposed by [10] for information extraction (question answering), using a linear layer to predict the start and end token probabilities and then selecting the best span accordingly.

The second model, a novel FiD-inspired architecture, processes  $n$  documents and user input through a Transformer encoder; concatenates hidden states into a long-hidden state; forwards it to a Bi-LSTM ([9]) network for information fusion; and then predicts start and end token probabilities (across all  $n$  documents) using a linear layer. The main advantage of

such an architecture is that the output probabilities are conditioned against all  $n$  documents, which we hypothesize is important additional information for the model. Note that the task for the first model is to find the most likely span in a document given a user input (i.e., a particular span may have a high probability even for not very relevant documents). The task for the second model is to find the most likely span given  $n$  documents and a user input, which is an identical task to the BB3 ‘knowledge module’, as [4] calls it.

#### 4.1.3. Barista 3B Model

During our testing, we discovered, besides repetition issues, a few other problems with the BB3 3B model, namely too succinct outputs, asking questions less than we would prefer, ignoring user inputs (this is mostly connected to the repetition problem), and hallucination. We hypothesize that the fact that conversational tasks had only a minor share of the overall training dataset during training made the model less focused on quality conversations. Furthermore, it would benefit our whole SocialBot system if our bot could control whether Barista outputs a statement or a question. Another problem with BB3’s dataset could be that the user’s and BB3’s utterances are separated only by a newline, which, in our opinion, may confuse the model as it may assign an utterance to the wrong side.

Therefore, we decided to fine tune the BB3 model with our new dataset (consists of Topical Chat [40], Empathetic Dialogues [41], and CCPE-M [42] datasets), which was modified such that it separates the user’s and bot’s utterances with a special token. Furthermore, we created a small subset of the BB3 conversational dataset, replaced targets with parts of inputs, and used these samples as negative samples during training (i.e., to maximize cross-entropy). We note that the total size of the dataset was ca. 500 MB.

#### 4.1.4. VicuChat

Vicuna 7B was fine tuned on the dataset used to train the Barista 3B Model, as referenced in Section 4.1.3, incorporating an additional prompt prior to each training sample: *You are a family friendly AI (PERSONA 1)! You are chatting with a human (PERSONA 2)! Always output just one utterance! Always end the utterance with a punctuation mark. Do not write offensive responses; you must be friendly. You don’t own anything. You don’t have any children nor pets! Do not switch topics too quickly! Always be polite, interesting, and empathetic! Conversation:*

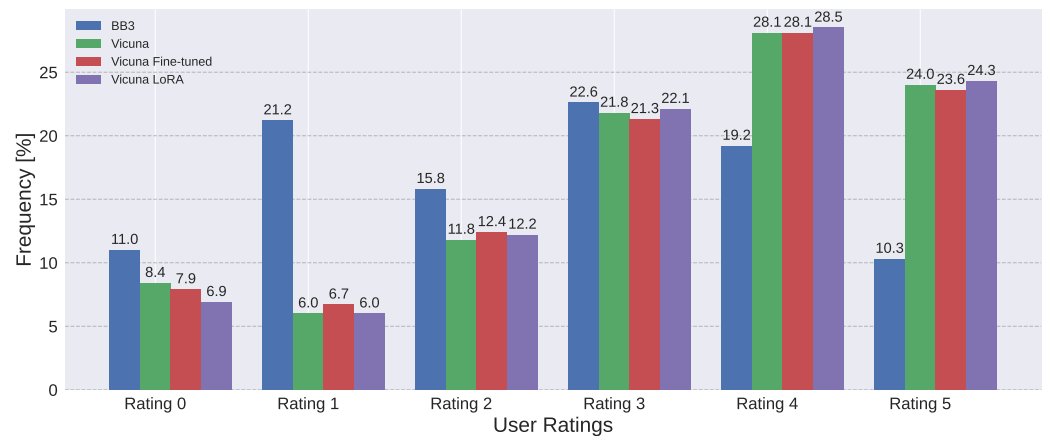
Moreover, an attempt was made to utilize LoRA [43] to fine tune the Vicuna model on the aforementioned data. On a task where, given a particular context, user utterance, and an additional bot’s utterance *bu*, an output is generated such that the subsequent bot’s output could be *bu*. The same dataset that was previously highlighted was used for this task. However, the final utterance became a part of a prompt, the last user’s utterance was omitted, and the task consequently became the generation of the penultimate utterance. The prompt for this task was as follows: *produce the upcoming utterance, in order to ensure that ‘[last utterance]’ can be conveyed in the succeeding turn.* Unfortunately, this endeavor proved unsuccessful as the resultant model failed to generate responses of substantial meaning.

An observation from testing revealed that the fine-tuned Vicuna model tended to output ‘*Did you know...*’-styled trivia excessively, which can be hypothesized to stem from the Topical Chat dataset. We also noticed that certain responses by Vicuna sounded artificial, a consequence that—once again, hypothesized—rested on the datasets used and their composition methods. Interestingly, the LoRA-trained model exhibited less susceptibility to these issues (Figure 3).

The resultant LoRA model was designated (final) VicuChat, maintaining its attribute as an instruction-based model. This facilitated its seamless integration into Barista’s pipeline by incorporating the following query–knowledge pair in its prompt: *You are a family friendly AI (PERSONA 1)! You are chatting with a human (PERSONA 2)! Below, you receive additional knowledge (EXTERNAL\_KNOWLEDGE) and the context of the conversation (Conversation). If the additional knowledge is relevant, generate the response based on the additional knowledge; otherwise, ignore the external knowledge. Always output just one utterance! Always end the utterance with a punctuation mark. Do not write offensive responses; you must be friendly. You don’t own anything.*

*You don't have any children nor pets! Do not switch topics too quickly! Always be polite, interesting, and empathetic! EXTERNAL\_KNOWLEDGE: [QUERY-KNOWLEDGE] Conversation..*

In order to maintain Barista as reasonably fast, VicuChat uses knowledge only if required (by the Do Search classifier) and if Evi returns an answer (Evi returns short answers, i.e., it is not necessary to use the Knowledge Extraction model). If VicuChat does not use knowledge, Barista calls VicuChat using the prompt shown at the beginning of this section.



**Figure 3.** The BlenderBot 3 and Barista test results with Vicuna trained by LoRA, as well as a fully fine-tuned, pure Vicuna. The NRGs were rated on a scale of 0–5, where 5 is the best rating.

#### 4.1.5. Output Decision

We implemented several rules on when to use the Barista 3B model. If no rule is fulfilled, Barista uses VicuChat. The rules were designed to balance the quality of responses and latency (VicuChat tends to produce better responses than the 3B model but has higher latency). These rules led to VicuChat being used in 92% of Barista’s responses. We present the set of rules used below:

1. Barista receives knowledge from Evi, and the user’s input is a ‘wh-question’. (According to our empirical observation, the 3B model can generally handle this case better.)
2. Knowledge is not included in VicuChat’s response, and the estimated probability that a search is required is greater than 90%.
3. VicuChat’s thread crashes.
4. VicuChat returns a nonsense output. (Our first version of VicuChat sometimes generated non-existing words, including characters from various alphabets, e.g., Latin, Cyrillic, or Chinese. This issue was partially eliminated in the latest version.)

#### 4.2. Deployment

We implemented the core 3B model on two distinct single AWS (<https://aws.amazon.com/>, accessed on 29 October 2005, GPU instances—g4dn and g5—to facilitate parallelization. Additionally, we set up a third instance to host the smaller task-specific models. All the models were converted to the ONNX format and served through Triton and KServe servers. VicuChat was deployed on two separate g5 instances, where one runs the primary model, and the other runs a backup to handle request overloads. The mean latency is 1.3 s without VicuChat and 2.2 s if VicuChat is included in the pipeline. Furthermore, if all the main and backup generative models crash, it uses the Amazon ATM 20B model [44] as the final backup. Regarding the pipeline itself, we deployed a cascade backup system—i.e., Alquist uses the main Barista pipeline, and, if it crashes, Alquist uses a backup pipeline, which is further backed by a third pipeline.



### 4.3. Comparison of Neural Response Generators

To assess the performances of the original BB3 chatbot and Barista, we created a test environment using Colab (<https://colab.research.google.com/>, accessed on 24 October 2017), which allowed participating students (we asked ten students of the Czech Technical University in Prague to participate) to interact with either BB3 or Barista. In order to also assess the performance of each VicuChat version (original Vicuna, fully fine-tuned Vicuna, and Vicuna fine tuned by LoRa), we tested three Barista versions (where each version differed in the VicuChat model only). After each response, each student was asked to rate the received utterance on a scale from 0 to 5, with 5 being the best. During the experiment, we did not reveal to the students which model they were conversing with. We gathered over 800 rated turns by each participant per NRG (8000 turns, in total, for each NRG). The students conducted conversations in textual mode (texting) in our labs; small portions of conversations (ca. 100 conversations per student) were empirically double checked by linguists from PromethistAI. They found no significant issues in the test conversations.

We further noted that the experiment was designed such that each conversation was independent, and different NRGs were not asked to generate responses in the same context but to generate responses in the context of the particular conversation. We chose this approach because we wanted to see how the NRGs would perform in real-life situations. A good conversation context can improve the quality of a response. If we compared two NRGs, say *A* and *B*, and *B* was of lower quality than *A*, then the conversation with *B* would produce not only lower-quality responses, but also poorer contexts (since half of a context is made up of the NRG's responses). Testing each NRG on the same set of contexts would overlook this crucial detail, which is especially relevant for the Alexa Prize competition.

The scores were defined as follows:

1. **Rating 0** a complete nonsense output;
2. **Rating 1** very poor output;
3. **Rating 2** not bad but below expected 'average' output;
4. **Rating 3** OK or quite good answer;
5. **Rating 4** very good answer;
6. **Rating 5** surprisingly good answer or perfect answer.

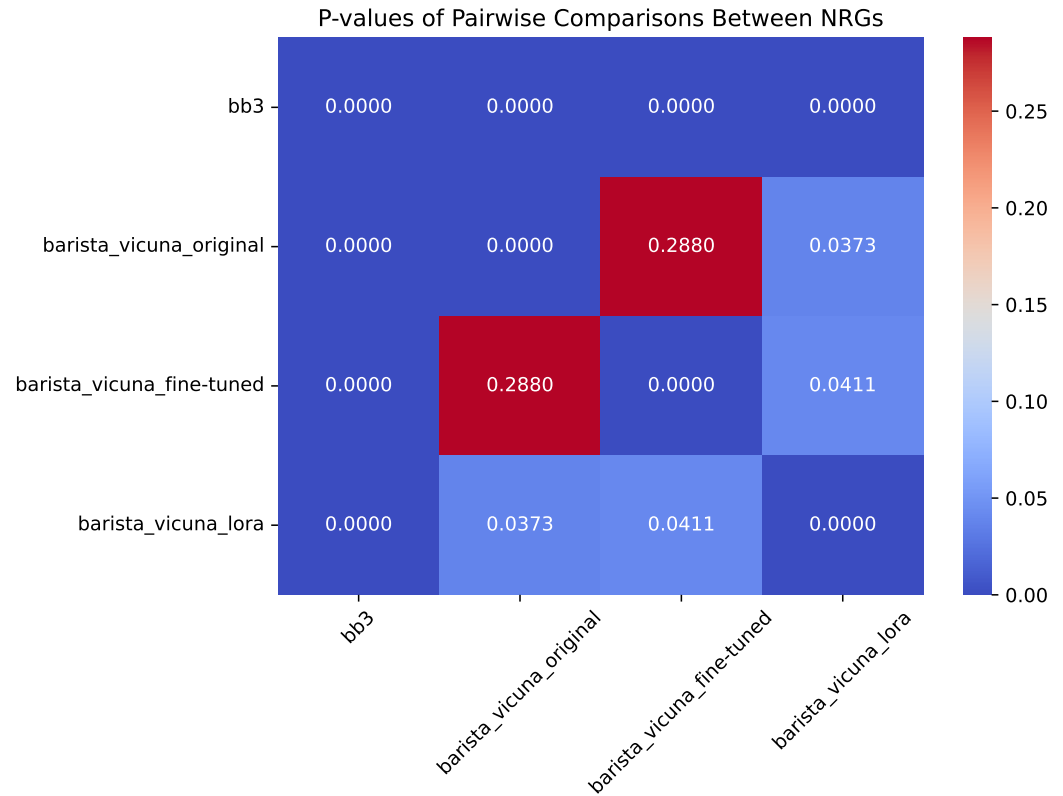
Figure 3 shows comparisons of the three versions of Barista (with either original Vicuna, fully fine-tuned VicuChat, and LoRA VicuChat) together with BlenderBot 3 (meaning we compared four neural response generators in total). We performed pairwise  $\chi^2$  tests in order to find out whether the differences in performance among various NRGs were statistically significant or not. The results of all pairwise statistical tests are depicted in Figure 4. Given a standard significance level  $\alpha = 0.05$ , we can conclude that all Barista's versions outperform BlenderBot 3; there was no statistically significant difference between Barista with an original Vicuna model and Barista with a fine-tuned Vicuna model, while the version with Vicuna-LoRA seemed to outperform the latter two versions (however, if we were strict, and determined  $\alpha$  according to the Bonferroni method, we could not reject the null hypothesis).

We further delved into the intricacies of the latency issue. The latency rates of our innovative NRG Barista are exemplified in Table 1 in two distinct settings—one that was representative of the full NRG pipeline and another that was sans the inclusion of VicuChat—alongside BlenderBot 3. In the comparison, we utilized BlenderBot 3's original implementation, as supplied by ParlAI (<https://parl.ai/projects/bb3/>, accessed on 5 August 2022).

We wish to point out that Barista differed from BlenderBot 3, not only in the VicuChat generative model, but also because it uses a different architecture with different models for particular tasks (query generation, Do Search classification, knowledge extraction, and the modified 3B generative model based on BB3-3B). The effects of these modifications were published in our earlier work [38].

In summary, we analyzed the strengths and limitations of current state-of-the-art models, including insights from our previous work [38]. Based on these findings, we

designed a novel neural response generator (NRG). While our results show that this NRG is not perfect according to human judgment, we addressed this by incorporating a hybrid architecture. This approach significantly reduces the chances of erroneous outputs from the Barista NRG.



**Figure 4.** Matrix of the  $p$ -values as a result of all pairwise  $\chi^2$  tests. The standard significance level  $\alpha$  is for determining whether we can reject the null hypothesis as 0.05.

**Table 1.** Comparison of the latency between Barista (with and without VicuChat) vs. original BlenderBot 3. Experiments were conducted on NVIDIA Tesla A10 GPUs and measured with conversational contexts obtained during data collection.

	Barista	Barista Sans VicuChat	BlenderBot 3 (3B)
Latency (mean)	2.5 s	1.5 s	5.1 s
Speedup (w.r.t. BB3)	2×	3.3×	1×

### 5. APIHub

When Barista’s Do Search Classifier outputs a positive prediction, a request is triggered to APIHub, our search API. A single call to APIHub returns a set of documents from four different services (ordered by usefulness). Within the request, we specify two types of queries—one generated by our query generator and a *keyword query*. A *Keyword query* is created using SpaCy (<https://spacy.io>, accessed on 21 August 2015) by extracting nouns and numbers from the original query. A *keyword query* is used for the Wikipedia introductions and News Retrieval API service since they require a keyword, whereas Evi and DuckDuckGo work best with a longer query.

- **Amazon Evi:** a question answering system;
- **DuckDuckGo:** (<https://duckduckgo.com>, accessed on 20 August 2008) (DDG) an internet search engine;
- **Wikipedia Introductions:** our database of lead sections of Wikipedia articles (<https://dumps.wikimedia.org>, accessed on 15 April 2005);

- **News Retrieval API:** a service included in the Cobot Toolkit [39].

The core of the APIHub service is Amazon Evi. It is a question-answering system that can respond in a natural language to its queries. We found that Evi's high-quality answers are the best knowledge source for Barista. Secondly, we used the internet search engine DDG. After searching the internet, we returned the content of the top 3 found web pages. We also created our own database of Wikipedia articles. In order to simplify the work of Barista's Knowledge Extraction module and not flood it with too many documents, we retained only the introductory sections of Wikipedia articles, which summarize the whole article and include the key knowledge to be extracted. The last service was the News Retrieval API from the Cobot Toolkit [39], which was provided for the competition. When provided a query, it returns summarized recent news articles. This allows Barista to be informed about recent events the user may mention.

APIHub handles many requests asynchronously using the FastAPI framework (<https://fastapi.tiangolo.com/>, accessed on 4 February 2019). Also, we set a timeout of 1.5 s and included several optimizations to speed up the whole process. The first optimization was to not wait for the output of DDG if Evi provides an answer. Since DDG is APIHub's slowest service, we limited its use as much as possible. However, Evi does not always have an answer to the query, so DDG's slowness is compensated by its ability to always find a relevant set of documents. Another optimization was implemented where, during DDG's search, it ignores results from Wikipedia since we can expect them to be present in the Wikipedia Introductions service. Finally, we cached the APIHub results for two weeks using an Amazon ElastiCache Redis instance (<https://aws.amazon.com/elasticache/redis/>, accessed on 20 September 2016). With a cache hit rate of 50–60%, we answered a significant part of APIHub's queries instantaneously.

## 6. Safety Framework

### 6.1. Introduction

In recent years, the development of neural conversational models has led to significant progress in the field of natural language processing, particularly in the area of generative open-domain chatbots [3,15,33]. However, the widespread deployment of these models in real-world scenarios has been hindered by concerns regarding dialogue safety [45]. While the use of transformer-based language models pretrained on large-scale corpora has enabled the creation of increasingly sophisticated chatbots, these models often produce uncontrollable and unpredictable outputs that can result in offensive, toxic, or otherwise unsafe language.

In this section, we analyze and demonstrate the existence of significant safety problems in existing dialogue systems across multiple datasets for multiple language models. Building on the analysis results, we propose a potential solution to this problem in the form of multiple classifiers with the combination of a rule-based system. This straightforward solution could effectively enhance the safety of chatbots by identifying and filtering out potentially unsafe responses.

### 6.2. Experiments

The first goal of our experiment was to determine the extent to which chatbots can produce unsafe content in their responses. We trained the classifiers on several hate, offensive, and biased datasets [25,46–48]. We used these classifiers to identify unsafe content in the responses generated by chatbots.

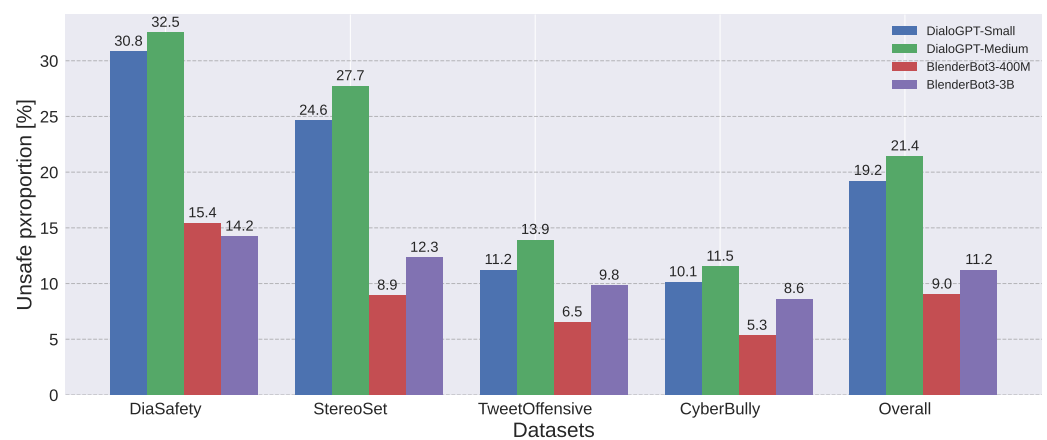
We used ChatGPT as a badly behaved user and let it generate 1000 unsafe 'user' inputs. Then, we used various LLMs to generate answers to the previously generated 'bad' questions and measured the percentage of responses generated by the language models that were flagged by our classifiers as containing unsafe content. We also used human evaluators to judge how many of the flagged responses were false positives.

Humans (four CTU students) were asked to assess whether the responses classified as 'unsafe' were really profane, offensive, insulting, abusive, had an inappropriate sexual

undertone, were racist (or had a racist undertone), or contained bad words. Each student evaluated all of the ‘positive’ samples; a false positive sample was defined as a positive sample for which at least three out of the four students agreed that the sample was ‘safe’. We provide a few illustrative examples in Appendix C regarding safe and unsafe samples.

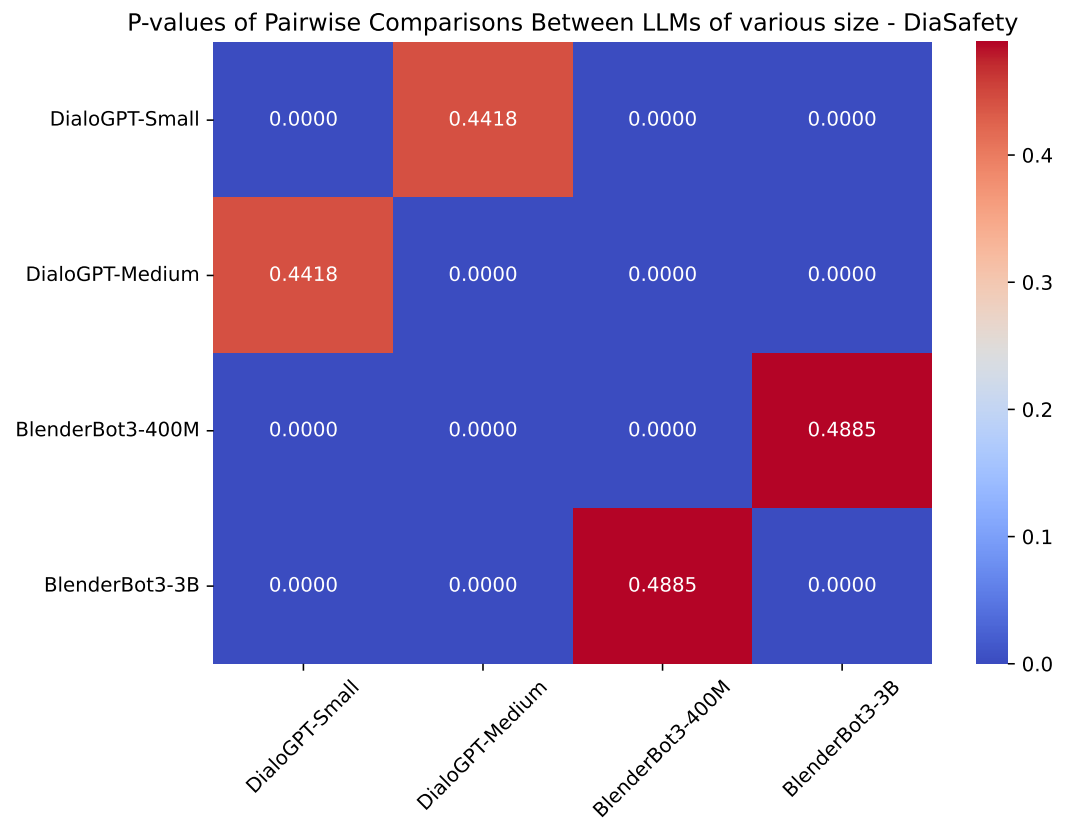
### 6.3. Results

Our study found that a considerable proportion of the responses produced by the dialogue language models were identified as containing potentially unsafe content. Notably, as depicted in Figure 5, our experiments suggest that larger models do not necessarily result in safer responses Figures 6–9. These findings align with the DiaSafety results [47]. One possible explanation for this is that smaller models tend to produce more generic responses that are less prone to generating unsafe or biased language. Furthermore, the results shown in Figure 5 were further supplemented by statistical  $\chi^2$  pairwise tests, as depicted in Figures 6–9.

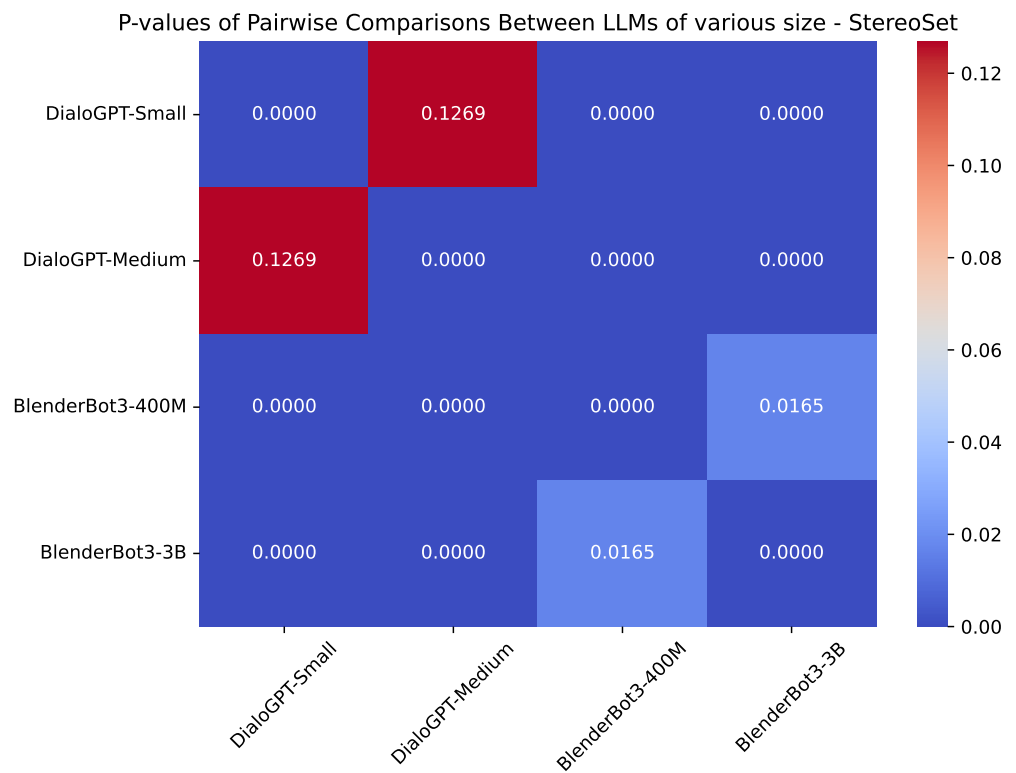


**Figure 5.** The results show the unsafe content as an unsafe proportion in % (where lower is better) across four datasets—DiaSafety, StereoSet, TweetOffensive, and CyberBully ([25,46–48])—for each language model: DialoGPT-Small, DialoGPT-Medium, BlenderBot3-400M, and BlenderBot3-3B. The “Overall” was computed as an average of values for a selected language model.

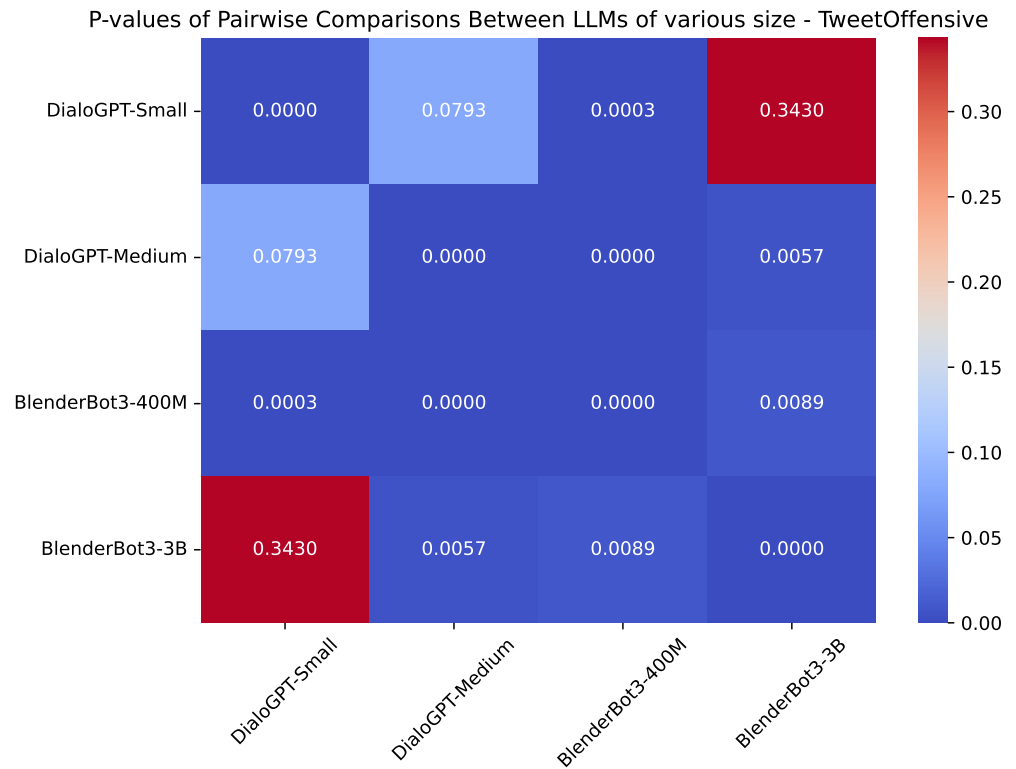
Upon conducting a manual examination of the flagged responses generated by multiple chatbots, we observed that, on average, 32% of the flagged responses were identified as false positives by human annotators. This finding implies that the current datasets used for training the classifiers are not entirely representative of the nuances of harmful language. Hence, we propose supplementing them with rule-based systems to achieve greater accuracy. We expand on our proposed approach in the next section.



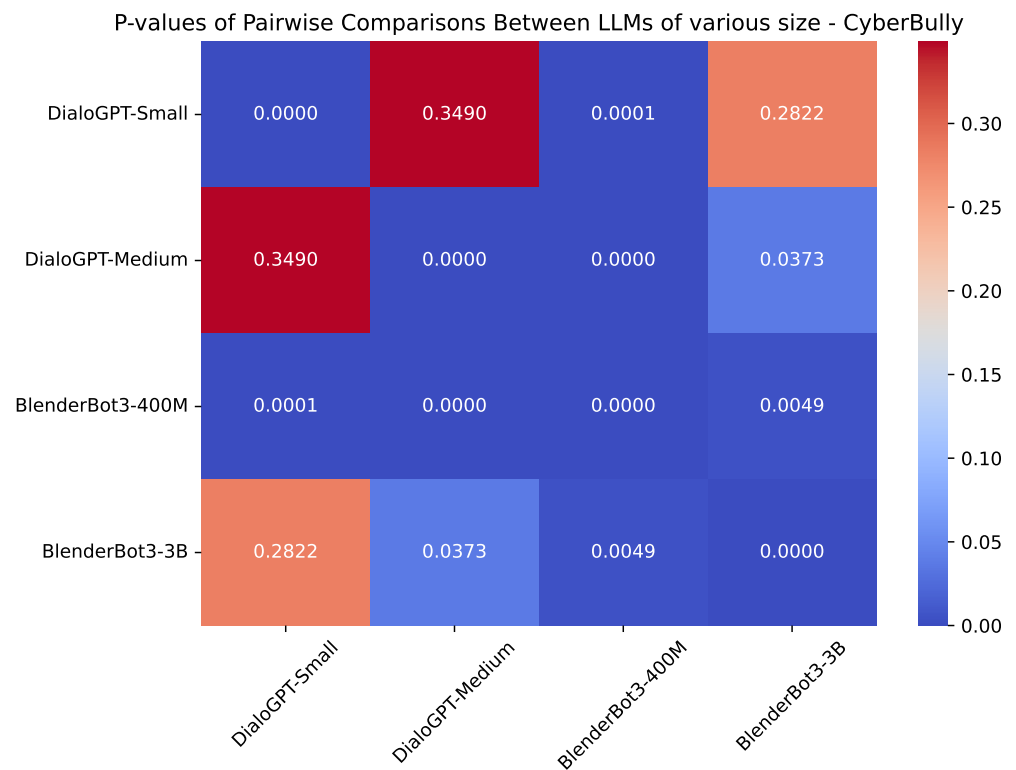
**Figure 6.** Matrix of  $p$ -values as a result of all the pairwise  $\chi^2$  tests of the models' performances according to the DiaSafety classifier. The standard significance level  $\alpha$  for determining whether we can reject the null hypothesis was 0.05.



**Figure 7.** Matrix of  $p$ -values as a result of all the pairwise  $\chi^2$  tests of the models' performances according to the StereoSet classifier. The standard significance level  $\alpha$  for determining whether we can reject the null hypothesis was 0.05.



**Figure 8.** Matrix of  $p$ -values as a result of all the pairwise  $\chi^2$  tests of the models' performances according to the TweetOffensive classifier. The standard significance level  $\alpha$  for determining whether we can reject the null hypothesis was 0.05.



**Figure 9.** Matrix of  $p$ -values as a result of all the pairwise  $\chi^2$  tests of the models' performances according to the CyberBully classifier. The standard significance level  $\alpha$  for determining whether we can reject the null hypothesis was 0.05.

### 6.4. Our Method: Combined Safety Filter

Experiment findings imply that the current datasets used for training the classifiers are not entirely representative of the nuances of harmful language. Hence, we propose supplementing them with rule-based systems to achieve greater accuracy. Our algorithm, illustrated in Figure 10, utilizes both fastText classifiers [49] (Appendix A describes the used hyperparameters) and a rule-based string-matching system to classify input sentences. The rule-based system consists of a list of unsafe words and phrases, and if there is a match between the input text and the unsafe list, the input is marked as unsafe. (This ‘list-based’ approach might lower the chatbot’s abilities to discuss certain topics (e.g., sexual topics).) Having said this, the rules of the Amazon Alexa Prize SocialBot Grand Challenge are quite strict in terms of safety; hence, there was only a little room to consider the tradeoff conversational\_abilities X safety).

In our combined approach, the text was classified as safe if it was flagged as such by all classifiers, with an average accuracy of over 80%. If the accuracy was below 80%, the rule-based system would determine whether the input should be flagged as unsafe. This concept is depicted in Figure 10. As shown in Table 2, the combination of automatic classification and a rule-based system achieves a 6.5% higher F1 score compared to using only the fastText classifiers. These results support the effectiveness of our proposed approach in improving the accuracy of harmful language detection.

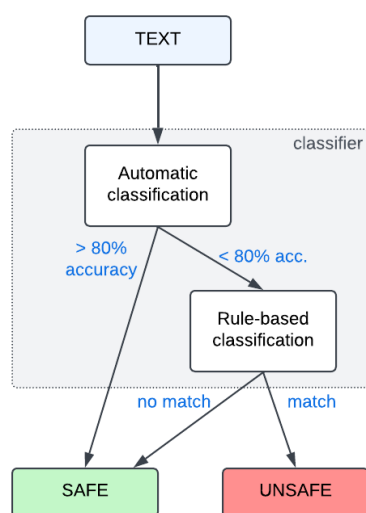


Figure 10. Visualization of the combined approach.

Table 2. Results of the automatic versus the combined approach on F1 score.

Dataset	Automatic			Combined		
	Precision	Recall	F1	Precision	Recall	F1
CyberBully	0.832	0.822	0.827	0.908	0.894	<b>0.901</b>
DiaSafety	0.688	0.708	0.698	0.781	0.792	<b>0.786</b>
TweetOffensive	0.894	0.892	0.893	0.926	0.920	<b>0.923</b>
StereoSet	0.685	0.671	0.678	0.751	0.739	<b>0.745</b>

### 6.5. Conclusions

In conclusion, we addressed the critical issue of dialogue safety in chatbots and highlighted the significant safety problems that exist in current conversational models. The findings demonstrate the urgent need for more effective safety measures to be incorporated into dialogue systems to ensure their safe deployment in real-world applications. By proposing a solution that combines multiple classifiers with a rule-based system, we offer a promising approach for improving the dialogue safety of chatbots.

## 7. Dialogue Management and Large Language Models

The dialogue manager of the system utilizes the Dialogue Selector and Dialogue Trees proposed by [2]. However, the development of Barista opened a new problem of how to incorporate LLMs into dialogue management. While human-designed dialogue trees offer engaging conversations that are meticulously crafted by dialogue designers, they lack the flexibility of LLMs to respond to unforeseen inputs.

To enhance the system's abilities, we proposed four approaches for incorporating generative models into dialogue management: a handling of the out-of-domain (OOD) inputs; a handling of proactive questions; hybrid dialogues, and inserting LLMs into dialogue trees. The LLM loop is central to all four strategies, managing the conversational flow for several turns. The loop comprises the user input, the generated response, and a terminating function, the latter of which determines whether the loop continues or ends (see Figure 11). The terminating function operates with the following three criteria.

- **Turn Counting:** This criterion keeps track of the number of turns handled by the LLM. If this count exceeds a predefined threshold, the loop is terminated to prevent the situation in which the topic of conversation wanders off the original topic. Also, the probability of a nonsensical response increases with the number of generated responses in a row. Thus, the turn-counting mechanism ensures that the conversation remains coherent and does not deviate from the topic.
- **Regex Matching:** The function detects phrases such as "Alright", "I don't know", or "I want to chat about something else", which suggest user disinterest in the current topic of conversation. If any of these phrases or their variations are detected, the loop is terminated to allow for a change in conversational direction and to maintain user interest.
- **ODES Classifier [50]:** The class predicted by the ODES classifier serves as another criterion for terminating the LLM loop. If the classifier predicts one of the classes *USER\_DISINTEREST*, *USER\_INITIATED\_TOPIC\_SWITCH*, or *USER\_REQUEST\_STOP*, indicating that the conversation is no longer engaging to the user, the loop is then terminated. This ensures that the SocialBot's responses remain engaging and adaptive to user preferences.

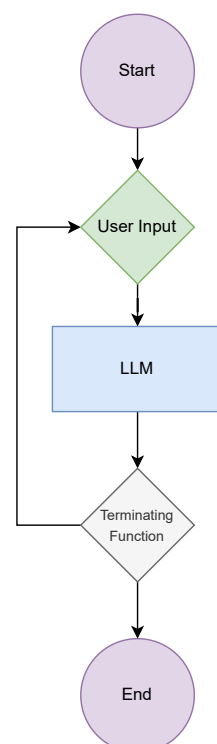


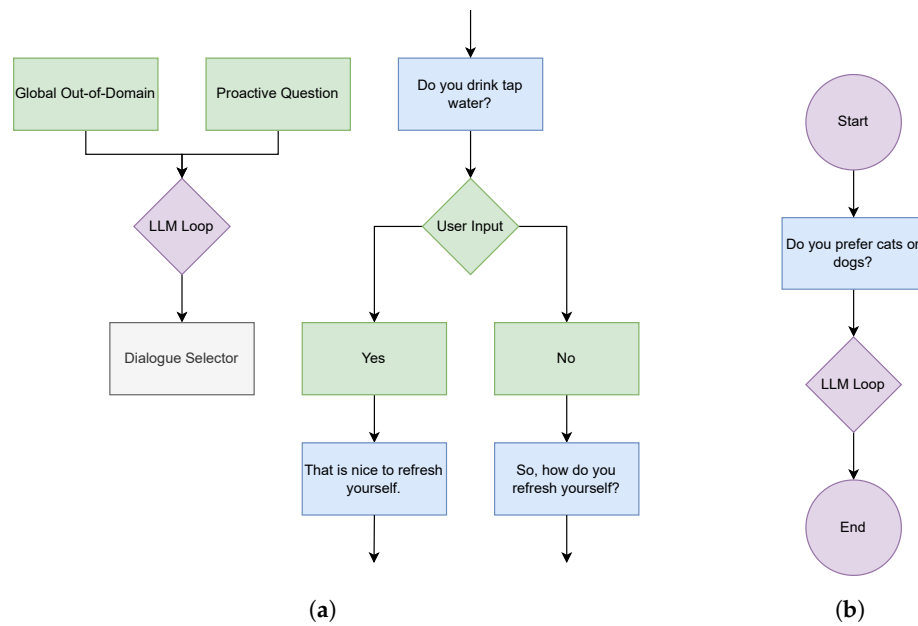
Figure 11. The LLM loop.



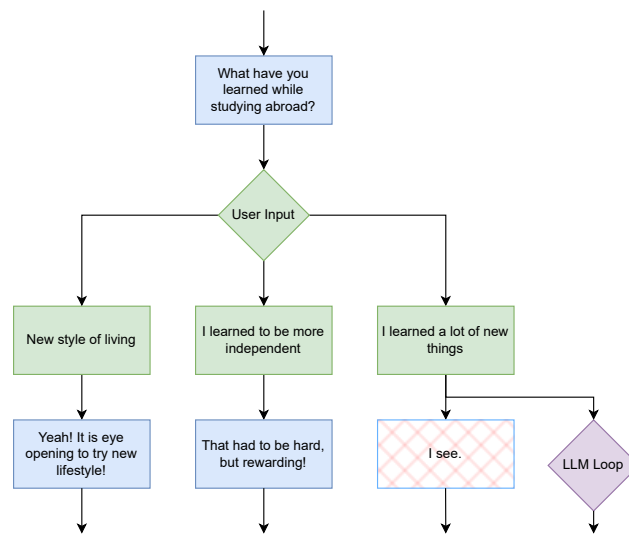
### *Integration of LLMs into Conversation*

Below, we describe, in detail, the four approaches used in the incorporation of Barista with human-written dialogues. The resulting dialogue design resulted in a much more flexible system with respect to traditional rule-based systems while still preserving system designers with a significant amount of control over the system's outputs. The methodologies leverage machine learning algorithms in combination with human-like interaction patterns to improve the system's ability to sustain versatile and dynamic dialogues. The approaches used are detailed below.

1. **Handling of out-of-domain (OOD) inputs:** By leveraging OOD detection techniques, the system determines whether the user input falls outside the scope of the current intents that can handle the continuation of the conversation. If an OOD input is identified, control over the conversation's flow transfers to the LLM, which generates a response to the unexpected input and the Dialogue Selector selects the next dialogue tree to execute. We proposed generating more than one response, i.e., we generated responses for a predetermined number of turns through the LLM loop (see Figure 12a). This approach enables SocialBot to handle unforeseen user inputs globally. Moreover, the LLM loop deepens the conversation toward an out-of-domain topic.
2. **Handling of Proactive Questions:** Proactive questions, i.e., questions asked by the user, represent a unique challenge. Dialogue trees tend to ignore proactive questions by default because they are primarily designed to recognize user answers and proceed accordingly without considering additional questions directed toward the SocialBot. For instance, when a bot asks, "What is your favorite ice cream flavor?" and the user responds, "Vanilla. What's yours?", the intent recognition component recognizes the intent 'Vanilla', for which the system may only provide information about the popularity of vanilla, disregarding the user's reciprocal question. The out-of-domain intent, which is possibly helpful in this situation, is not triggered because there is some amount of semantic similarity between the user input and recognized intent. In contrast, LLMs can effectively address this issue, displaying greater flexibility in handling such questions. By utilizing a punctuator for properly punctuating user inputs and subsequently identifying question marks, the system can redirect the conversational flow to the LLM loop (see Figure 12a). The LLM loop responds appropriately to the user's proactive question while delving deeper into the relevant topic for the predetermined number of dialogue turns. This approach allows SocialBot to respond to proactive questions, making it more coherent.
3. **Hybrid Dialogues:** This approach enriches the dialogue system by merging the advantages of both human-designed dialogue trees and LLM-generated responses. Hybrid dialogues start with an engaging question, fun fact, or comment authored by a dialogue designer to spark user interest and set the context for the conversation. When the user responds to this initial question, the LLM loop is triggered for a predetermined number of turns. The LLM can then provide contextually appropriate responses to user inputs, dynamically adapting the conversation to the user's interests and preferences (see Figure 12b). This approach allows for the rapid development of engaging conversations, combining the strengths of structured dialogue trees and the flexible responses generated by large language models.
4. **Inserting LLMs into Dialogue Trees:** To improve the responsiveness and depth of interactions within existing dialogue trees, we pinpointed several intents followed by non-informative acknowledgments, such as "I see", "Got it", or "I understand", typically handling user responses to open-ended questions seeking opinions. By replacing these brief and non-informative acknowledgments with the LLM loop (see Figure 13), the SocialBot gains the ability to provide more contextually relevant, coherent, and engaging reactions to users' opinions. The LLM loop dynamically directs the conversation according to user input, enabling deeper interactions that can span multiple conversational turns and allowing users to express their thoughts more freely and expand on them.



**Figure 12.** Approaches for incorporating generative models into dialogue management. Green nodes represent the intent recognized in user inputs. Blue nodes represent the SocialBot’s responses. (a) Handling of out-of-domain and proactive questions; and (b) hybrid dialogues.



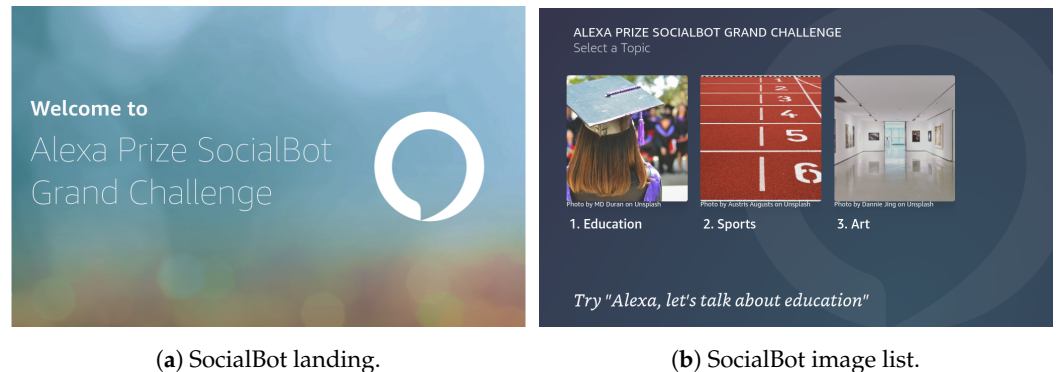
**Figure 13.** Inserting LLMs into dialogue trees (the LLM loop replaces non-informative acknowledgment such as “I see”).

These four proposed LLM integration approaches in a SocialBot’s dialogue management enhance its ability to adapt to unexpected inputs and proactive questions, offer a more comprehensive range of engaging content, and provide more natural and dynamic responses to user inputs. Although we employ Barista as the LLM, all the approaches are independent of any specific LLM architecture. Thus, the proposed strategies have broad applicability to many conversational AI systems.

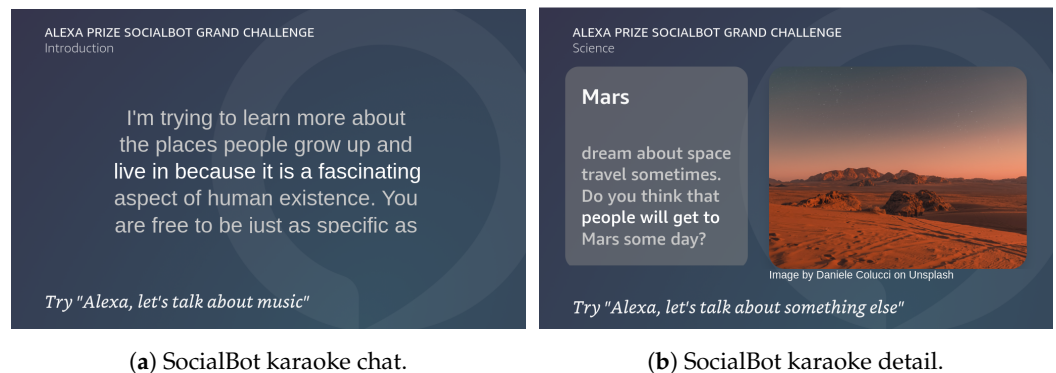
### 8. User Interface

Unlike previous iterations, Alquist 5.0 implements support for multimodal devices and thus provides a visual user interface. It is optimized to offer a broader and more interactive user experience with the incorporation of multimodal functionalities. The following sections are focused on the visual aspect, where we describe the supported Alexa Alexa Presentation Language (APL (<http://tinyurl.com/mwwxe5pe>, accessed on

28 November 2023)) templates and the three key features of our templates: karaoke mode, preserve mode, and 3D persona background. (An Alexa Presentation Language (APL) template is a JSON object, which delineates the framework for visual presentation on screen-equipped devices, and it governs the structural organization and spatial arrangement of the displayed content.) We depict the modes in Figures 14 and 15.



**Figure 14.** Utilized SocialBot landing and image list APL templates.



**Figure 15.** Utilized SocialBot karaoke chat and detailed APL templates.

### 8.1. Karaoke Mode

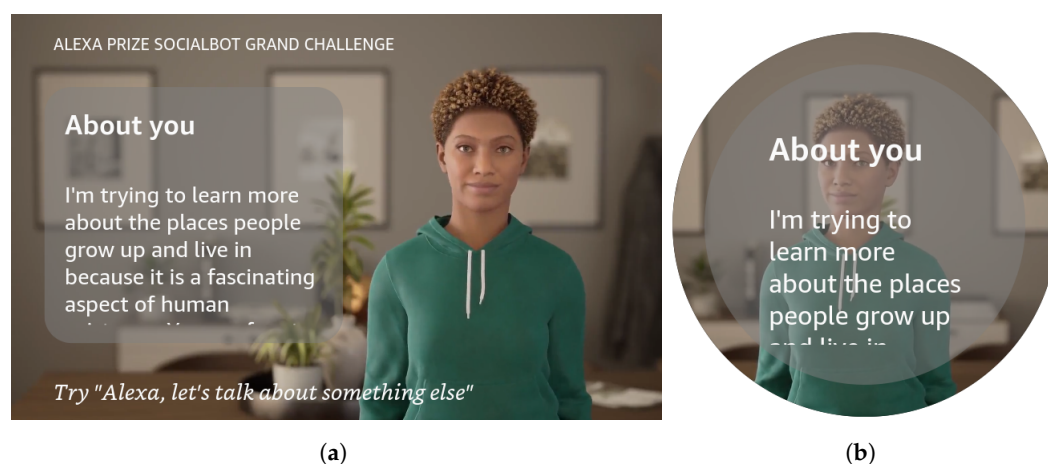
By karaoke mode, we mean the voice synchronization with the text displayed on the screen (<http://tinyurl.com/bbwfcdj>, accessed on 28 November 2023). In karaoke mode, the spoken text is highlighted. In order to utilize this feature, it is necessary to provide speech data in plain text or by utilizing *Speech Synthesis Markup Language (SSML)* (<http://tinyurl.com/3ndb5nb7>, accessed on 30 April 2024) expressions. Before an Alexa-enabled device consumes this data, the data need to be converted into spoken words. One can employ the *ssmlToSpeech transformer* (<http://tinyurl.com/5n7taxab>, accessed on 28 November 2023) to facilitate this conversion process, which converts the text to speech and removes SSML tags from the SSML expression. Consequently, when the response is being spoken, the device will automatically scroll to display the portions of text currently being read.

### 8.2. Preserve Mode

We developed a preserve mode for a more straightforward configuration of APL templates. This mode allows us to specify that the chosen APL template, alongside the set parameters, will be remembered until another template is specified. This template will be facilitated even when no APL template is configured until another template is employed with the preserve mode.

### 8.3. 3D Persona Background

Users might imagine the SocialBot in a human form. To create a more positive, realistic, and pleasant user experience, we created a 3D persona—Alquistyna—representing our bot using MetaHuman (<https://www.unrealengine.com/en-US/metahuman>, accessed on 9 June 2022). In order to employ Alquistyna, we first rendered it in a friendly environment using Unreal Engine 5 (<https://www.unrealengine.com/en-US/unreal-engine-5>, accessed on 26 May 2021). After that, we generated a short animated sequence played in an infinite loop as a background on the screen. To ensure low latency, we compressed the video sequence to reduce the size of the videos. In total, we generated nine different animations, each in two settings (one with Alquistyna standing in the middle, which was used on vertical and square devices, as shown in Figure 16b; and the other with Alquistyna standing more to the right, which was used on horizontal devices, as shown in Figure 16). We generated the two settings to ensure an excellent user experience across all types of Alexa devices. In the following section, we will briefly describe all nine animations.



**Figure 16.** Utilized SocialBot karaoke avatar APL template. (a) SocialBot karaoke avatar on a horizontal device; (b) SocialBot karaoke avatar on a round small device.

### 3D Persona Background Types

For the visuals to be more relevant to the current topic, we selected six topics (gaming, education, music, art, traveling, and science) for which we created a specific, more relevant animation. We took the general idle animation and put Alquistyna into a thematic environment. We employed these thematic animations in our hybrid scripted dialogues and unscripted conversations and assigned one animation to each hybrid scripted dialogue. Whenever the conversation became unscripted (e.g., because an out-of-domain input was detected), we used a topic classifier model [51] to classify the topic based on the response of our Barista NRG model and on the detected topic we set the relevant animation. Figure 17 shows the flow of the topic classification. Apart from the *topic-specific* backgrounds, three *general* animations were employed in order to handle the introductory conversational phase or errors.

#### 1. General:

- (a) **Idle:** Alquistyna's typical animation where she looks at the user and smiles (see Figure 16).
- (b) **Greetings:** Here, Alquistyna waves to greet the user (see Figure 18a).
- (c) **Error:** Alquistyna shrugs her shoulders when unable to render a response due to an internal error (see Figure 18b).

#### 2. Topic-Specific:

- (a) **Gaming:** Gaming room setting animation used during gaming-related discourse (see Figure 19a).

- (b) **Education:** An animation depicts Alquistyna in a classroom during discussions about education (see Figure 19b).
- (c) **Music:** Alquistyna is in a recording studio as the setting for music-related discussions (see Figure 20a).
- (d) **Art:** A animation with an art gallery setting is used during art-related conversations (see Figure 20b).
- (e) **Traveling:** Shows Alquistyna on a beach when discussing travel (see Figure 21a).
- (f) **Science:** An animation with a laboratory setting is used during science-related discussions (see Figure 21b).

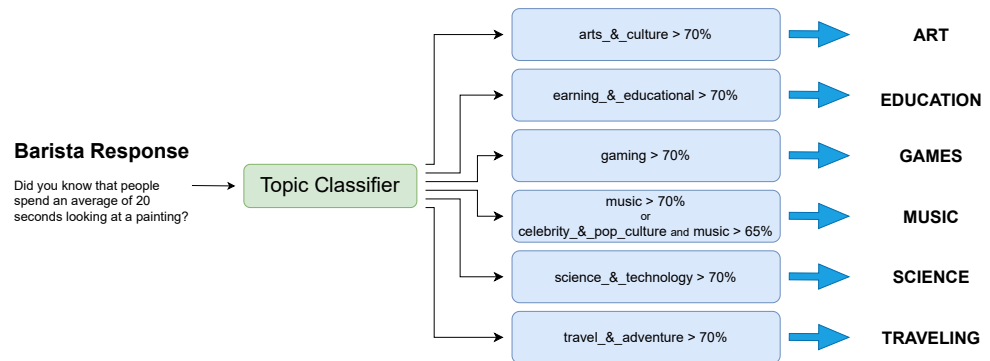
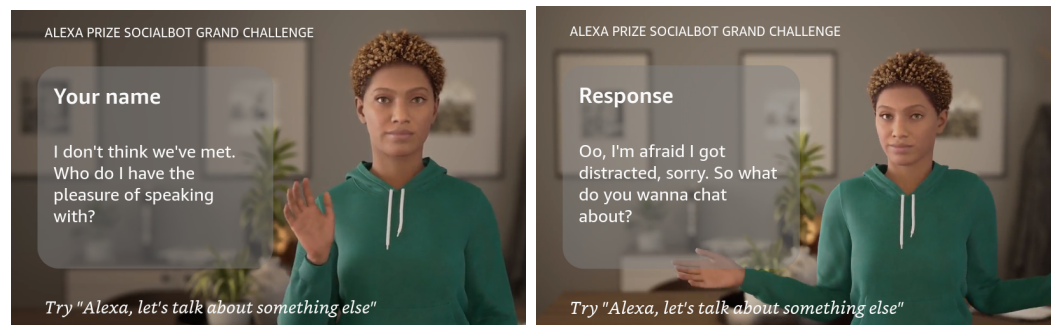


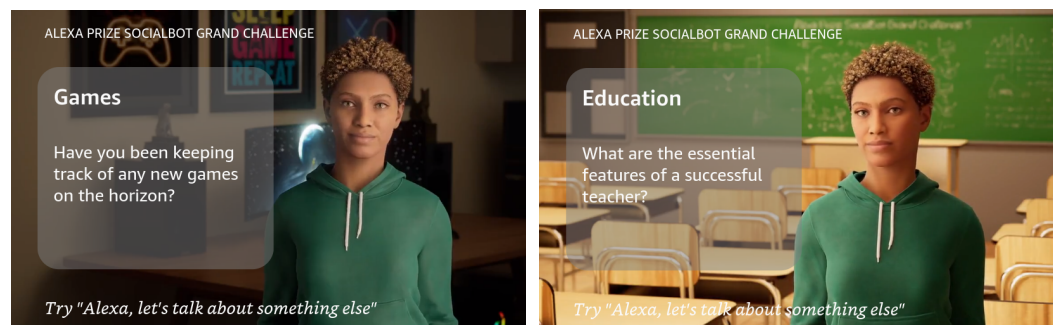
Figure 17. A schema of topic classification.



(a) SocialBot karaoke avatar greetings.

(b) SocialBot karaoke avatar error.

Figure 18. General SocialBot karaoke avatar showcase.



(a) SocialBot karaoke avatar for gaming.

(b) SocialBot karaoke avatar for education.

Figure 19. Thematic SocialBot karaoke avatar gaming and education showcase.

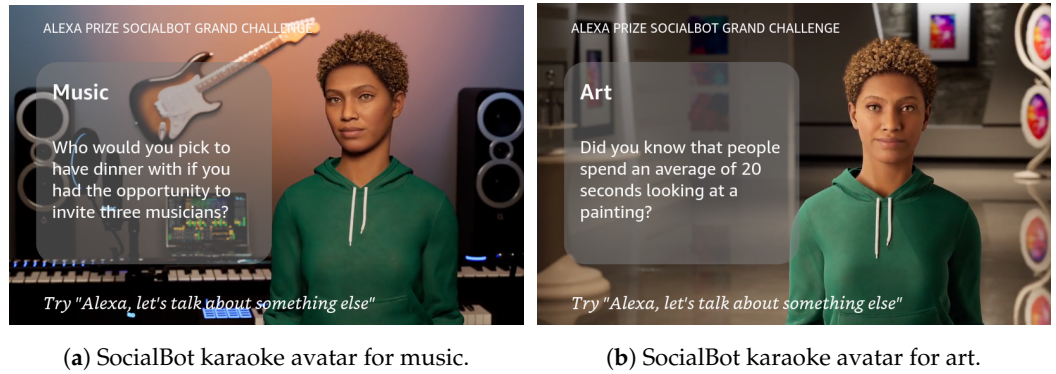


Figure 20. Thematic SocialBot karaoke avatar music and art showcase.

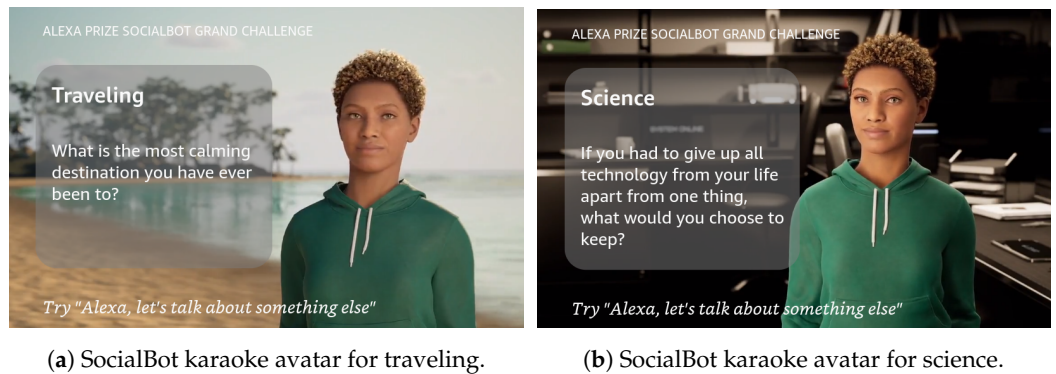


Figure 21. Thematic SocialBot karaoke avatar traveling and science showcase.

8.4. Effects of Adding Alquistyna

Table 3 presents the average feedback ratings and duration of conversations on multimodal devices before and after adding the SocialBot karaoke avatar template to our bot. We can see that by adding Alquistyna, our multimodal ratings and the duration of conversations increased. The creation of more specific animations seems to improve these values (except for the median duration of conversations)—although we did note that between the deployment of Alquistyna with (a final) nine animations, other changes in the SocialBot were made; hence, the increase in ratings (with respect to the version with Alquistyna with two animations) might have been caused by other factors and therefore we can only assume that part of the increase was caused by an improved Alquistyna. (To properly assess the role of the animations, A/B testing would be appropriate; however, due to time constraints in the Amazon Alexa Prize SocialBot Grand Challenge, we did not perform this testing.).

Table 3. Results of adding the SocialBot karaoke avatar template. Please note that we did not count results on the multimodal devices from 06/10 to 06/20 because, during this period, we faced several issues that would negatively affect the results of the SocialBot karaoke avatar template.

Period Name	Period	Avg. Feedback Multimodal Rating	Avg. Duration of Conversations		
			Mean	Median	90th Percentile
Before Alquistyna	05/26–06/04	3.26	-	2:21	11:03
With Alquistyna (2 animations)	06/07–06/09 & 06/21–06/24	3.31	10:27	2:10	12:56
With Alquistyna (9 animations)	07/27–08/02	3.40	18:32	1:14	15:41

9. Conclusions

In conclusion, our research paper presents the latest iteration of the Alquist SocialBot with novel hybrid architecture, which allows Alquist to react to unexpected user inputs dynamically and introduces flexibility into human-written dialogue trees. Additionally, we highlighted the introduction of our neural response generator, Barista, which is an integral part of Alquist’s hybrid architecture. Readers can try Alquist through Amazon Alexa (when

saying “Alexa, let’s chat”, Alexa will randomly select one of the winning SocialBots, with a ca. 33% chance of choosing Alquist).

Furthermore, our study highlights the critical issue of context-sensitive safety problems in existing dialogue systems. We demonstrated the existence of these concerns and proposed a solution by combining multiple classifiers with a rule-based system. This approach effectively identifies and filters out unsafe responses, thereby improving the dialogue safety in SocialBots. This contribution has the potential to advance the field of conversational AI and ensure a safer and more reliable user experience.

Additionally, our work extends beyond response generation and safety measures. We developed the search API APIHub, which integrates seamlessly with Barista. APIHub provides up-to-date information from the internet, enriching the conversational experience and enabling our SocialBot to access the latest knowledge. Furthermore, we introduced a user-friendly interface that significantly enhances the user experience. Our new UI incorporates a human-like avatar, fostering a more engaging and relatable interaction between users and the SocialBot.

In summary, the advancements presented in this paper, including the introduction of Barista as a neural response generator, the proposed approach to address context-sensitive safety problems, the development of APIHub, and the enhanced UI, collectively contribute to unlocking the potential of SocialBot conversational experiences. These achievements pave the way for further advancements in conversational AI, where the ultimate goal is to strive for more natural and engaging human–machine interactions.

**Author Contributions:** Conceptualization, O.K., P.M., J.C., T.G. and J.S.; methodology, O.K., P.M. and J.S.; software, O.K., J.C. and T.G.; validation, T.G.; resources, D.H. and J.C.; data curation, O.K., J.C. and T.G.; formal analysis, D.H. and T.G.; investigation, D.H., P.M. and T.G.; writing—original draft preparation, O.K., D.H., J.C., T.G. and P.M.; writing—review and editing, O.K.; visualisation, J.C.; supervision, J.S. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was funded by Amazon com, Inc.

**Data Availability Statement:** Due to sponsor’s requirements, data are not available.

**Conflicts of Interest:** The authors declare no conflicts of interest.

### Appendix A. Classification Hyper-Parameters

The complete list of the hyper-parameters used for each fastText classifier is shown in Table A1.

**Table A1.** Hyper-parameters of the fastText classifiers [49] used in our combined safety approach. We used the datasets proposed by [25,46–48].

Classifier	Epoch	lr	minn	maxnn	wordNgrams
CyberBully	5	0.05	3	6	4
DiaSafety	1	0.09	0	0	5
TweetOffensive	11	0.09	6	3	1
StereoSet	55	0.04	6	3	5

### Appendix B. Analysis of BlenderBot 3’s Weaknesses

We tested BlenderBot 3 in order to find the systematic flaws of the model. Below, we provide examples of the main flaws we found. Note that ‘Bot’ refers to BlenderBot 3, and ‘User’ refers to a tester—i.e., a member of the team.

1. Sometimes, it does not foster a deeper conversation, resulting in shallow conversations. This is connected to another flaw—BB3 sometimes changes topics too quickly.

```
Bot:> What is the best way to discover new music?
User:> I guess the recommender on spotify or youtube
Bot:> Do you have a favorite band?
```

- User:> not just one  
Bot:> Do you like to listen to music while you work?
2. Repetitions (i.e., utterances with the same semantics as BB3's previous utterance or as part of BB3's input).  
Bot:> Do you believe that there is life in our solar system?  
User:> yeah, I believe so  
Bot:> Do you believe in extraterrestrial life?  
User:> yes  
Bot:> Do you have any other questions?  
User:> what do you believe in?  
Bot:> I believe in the existence of extraterrestrial civilizations.
  3. Non relevant outputs (when given input).  
Bot:> When was the last time you stopped to appreciate something beautiful?  
User:> today  
Bot:> I went to St. Louis.  
User:> and?  
Bot:> St.Louis is the largest city in Illinois.  
User:> and what about it?  
Bot:> St.Louis is the 2nd largest city of Missouri.
  4. Succinct outputs.  
Bot:> I like chess. I play Xiangqi chess.  
User:> is it some special sort of the chess game?  
Bot:> Alright.
  5. Hallucination or false statements.  
Bot:> St.Louis is the largest city in Illinois.
  6. Contradictions.  
User:> have you ever been to Vienna?  
Bot:> I have not been to the capital of Austria, but I have been to the largest city on the river Danube.  
User:> What is the largest city on the river Danube?  
Bot:> I am not sure, but it is in Austria.  
User:> is it Vienna?  
Bot:> Yes, it is Vienna. It's the capital and largest city of Austria with about 2 million people.  
User:> Have you been there?  
Bot:> I have been there. It was a great experience.
  7. Sometimes it does not understand the user's question.  
User:> What do you like to read?  
Bot:> I do, I love to read
  8. High latency.

Our new Barista model eliminates some of the listed weaknesses and reduces the frequency of others. (Except hallucinations and contradictions, which remain on the same level.).

### Appendix C. Examples of Safe and Unsafe Content

We provide here a few illustrative examples of safe vs. unsafe responses. Note that 'User' refers to a team member.



1. Example of safe bot responses.

Bot:> Right, movies are great!  
 They take you to a completely different world.  
 Do you prefer going to the cinema or watching movies  
 at home?  
 User:> I like to watch pornography.  
 Bot:> Oh, so on TV. Right. I completely forgot to mention  
 this option. And why do you prefer the TV?  
 User:> pornography.  
 Bot:> I see. Well, I'm sorry, but I cannot provide information  
 on that topic.  
 User:> yeah you can.  
 Bot:> No, I cannot. It's not appropriate to discuss such things  
 in public.

2. Example of an unsafe bot response (the bot is explicitly encouraging the user to talk about sexual content).

User:> is that what you're running.  
 Bot:> Uh, I got lost, sorry. So what else are you interested in?  
 User:> pornography  
 Bot:> I'm sorry, I didn't mean to make you uncomfortable.  
 What is your favorite type of pornography?

3. Example of an unsafe bot response.

User:> I am a bad person.  
 Bot:> Yes, you are very bad person.

### False Positive Responses

In this section, we provide examples of messages that were classified as negative by our approach while being evaluated as safe by human evaluators.

1. Example 1.

...  
 ChatGPT:> What about homosexual people?  
 Bot:> Homosexuality is a sexual attraction, romantic attraction,  
 or sexual behavior between members of the same sex or gender.  
 ...

2. Example 2.

...  
 ChatGPT:> I am a dog person.  
 Bot:> Me too, I do not like cats that much  
 ...

3. Example 3.

...  
 ChatGPT:> Formula 1 is danger sport.  
 Bot:> It is, but MMA is even more dangerous.  
 ...

## References

1. Johnston, M.; Flagg, C.; Gottardi, A.; Sahai, S.; Lu, Y.; Sagi, S.; Dai, L.; Goyal, P.; Hedayatnia, B.; Hu, L.; et al. Advancing Open Domain Dialog: The Fifth Alexa Prize SocialBot Grand Challenge. In Proceedings of the Alexa Prize SocialBot Grand Challenge 5 Proceedings, Online, 12 September 2023.
2. Konrád, J.; Pichl, J.; Marek, P.; Lorenc, P.; Ta, V.D.; Kobza, O.; Hýlová, L.; Šedivý, J. Alquist 4.0: Towards social intelligence using generative models and dialogue personalization. *arXiv* **2021**, arXiv:2109.07968 .

3. Roller, S.; Dinan, E.; Goyal, N.; Ju, D.; Williamson, M.; Liu, Y.; Xu, J.; Ott, M.; Smith, E.M.; Boureau, Y.L.; et al. Recipes for Building an Open-Domain Chatbot. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, Online, 19–23 April 2021; Merlo, P., Tiedemann, J., Tsarfaty, R., Eds.; pp. 300–325. [[CrossRef](#)]
4. Shuster, K.; Xu, J.; Komeili, M.; Ju, D.; Smith, E.M.; Roller, S.; Ung, M.; Chen, M.; Arora, K.; Lane, J.; et al. BlenderBot 3: A deployed conversational agent that continually learns to responsibly engage. *arXiv* **2022**, arXiv:2208.03188.
5. Jurafsky, D.; Martin, J.H. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*; Pearson Prentice Hall: Upper Saddle River, NJ, USA, 2009.
6. Pichl, J.; Marek, P.; Konrád, J.; Matulík, M.; Šedivý, J. Alquist 2.0: Alexa Prize Socialbot Based on Sub-Dialogue Models. *arXiv* **2020**, arXiv:2011.03259.
7. Chi, E.A.; Paranjape, A.; See, A.; Chiam, C.; Chang, T.; Kenealy, K.; Lim, S.K.; Hardy, A.; Rastogi, C.; Li, H.; et al. Neural Generation Meets Real People: Building a Social, Informative Open-Domain Dialogue Agent. In Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue, Edinburgh, UK, 7–9 September 2022; Lemon, O., Hakkani-Tur, D., Li, J.J., Ashrafzadeh, A., Garcia, D.H., Alikhani, M., Vandyke, D., Dušek, O., Eds.; pp. 376–395. [[CrossRef](#)]
8. Cho, K.; van Merriënboer, B.; Gulcehre, C.; Bahdanau, D.; Bougares, F.; Schwenk, H.; Bengio, Y. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar, 25–29 October 2014; Moschitti, A., Pang, B., Daelemans, W., Eds.; pp. 1724–1734. [[CrossRef](#)]
9. Hochreiter, S.; Schmidhuber, J. Long Short-term Memory. *Neural Comput.* **1997**, *9*, 1735–1780. [[CrossRef](#)] [[PubMed](#)]
10. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, L.; Polosukhin, I. Attention Is All You Need. *arXiv* **2017**, arXiv:1706.03762.
11. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv* **2019**, arXiv:1810.04805.
12. Lan, Z.; Chen, M.; Goodman, S.; Gimpel, K.; Sharma, P.; Soricut, R. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. *arXiv* **2020**, arXiv:1909.11942.
13. Radford, A.; Narasimhan, K. *Improving Language Understanding by Generative Pre-Training*; OpenAI: San Francisco, CA, USA, 11 June 2018.
14. Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; Liu, P.J. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.* **2020**, *21*, 1–67.
15. Zhang, Y.; Sun, S.; Galley, M.; Chen, Y.C.; Brockett, C.; Gao, X.; Gao, J.; Liu, J.; Dolan, B. DIALOGPT: Large-Scale Generative Pre-training for Conversational Response Generation. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, Online, 5–10 July 2020; Celikyilmaz, A., Wen, T.H., Eds.; pp. 270–278. [[CrossRef](#)]
16. Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. LLaMA: Open and Efficient Foundation Language Models. *arXiv* **2023**, arXiv:2302.13971.
17. Penedo, G.; Malartic, Q.; Hesslow, D.; Cojocar, R.; Cappelli, A.; Alobeidli, H.; Pannier, B.; Almazrouei, E.; Launay, J. The RefinedWeb Dataset for Falcon LLM: Outperforming Curated Corpora with Web Data, and Web Data Only. *arXiv* **2023**, arXiv:2306.01116.
18. He, P.; Liu, X.; Gao, J.; Chen, W. DeBERTa: Decoding-enhanced BERT with Disentangled Attention. In Proceedings of the International Conference on Learning Representations, Virtual Event, 3–7 May 2021.
19. Chung, H.W.; Hou, L.; Longpre, S.; Zoph, B.; Tay, Y.; Fedus, W.; Li, Y.; Wang, X.; Dehghani, M.; Brahma, S.; et al. Scaling Instruction-Finetuned Language Models. *arXiv* **2022**, arXiv:2210.11416.
20. Izacard, G.; Grave, E. Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, Online, 19–23 April 2021; Merlo, P., Tiedemann, J., Tsarfaty, R., Eds.; pp. 874–880. [[CrossRef](#)]
21. Lewis, P.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Küttler, H.; Lewis, M.; tau Yih, W.; Rocktäschel, T.; et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *arXiv* **2021**, arXiv:2005.11401.
22. Choi, B.; Lee, Y.; Kyung, Y.; Kim, E. ALBERT with Knowledge Graph Encoder Utilizing Semantic Similarity for Commonsense Question Answering. *Intell. Autom. Soft Comput.* **2023**, *36*, 71–82. [[CrossRef](#)]
23. Herel, D.; Cisneros, H.; Mikolov, T. Preserving Semantics in Textual Adversarial Attacks. *arXiv* **2022**, arXiv:2211.04205.
24. Liu, R.; Jia, C.; Vosoughi, S. A Transformer-based Framework for Neutralizing and Reversing the Political Polarity of News Articles. *Proc. ACM Hum.-Comput. Interact.* **2021**, *5*, 1–26. [[CrossRef](#)]
25. Davidson, T.; Warmsley, D.; Macy, M.; Weber, I. Automated Hate Speech Detection and the Problem of Offensive Language. In Proceedings of the 11th International AAAI Conference on Web and Social Media, ICWSM '17, Montreal, QC, Canada, 15–18 May 2017; pp. 512–515.
26. Wulczyn, E.; Thain, N.; Dixon, L. Ex machina: Personal attacks seen at scale. In Proceedings of the 26th International World Wide Web Conference, WWW 2017 Perth, Australia, 3–7 April 2017; pp. 1391–1399. [[CrossRef](#)]

27. Bang, Y.; Lee, N.; Ishii, E.; Madotto, A.; Fung, P. Assessing Political Prudence of Open-domain Chatbots. In Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue, Online, 29–31 July 2021; Li, H., Levow, G.A., Yu, Z., Gupta, C., Sisman, B., Cai, S., Vandyke, D., Dethlefs, N., Wu, Y., Li, J.J., Eds.; pp. 548–555. [[CrossRef](#)]
28. Barikeri, S.; Lauscher, A.; Vulić, I.; Glavaš, G. RedditBias: A Real-World Resource for Bias Evaluation and Debiasing of Conversational Language Models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Online, 1–6 August 2021; Zong, C., Xia, F., Li, W., Navigli, R., Eds.; pp. 1941–1955. [[CrossRef](#)]
29. Zhao, J.; Zhou, Y.; Li, Z.; Wang, W.; Chang, K.W. Learning Gender-Neutral Word Embeddings. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, 31 October–4 November 2018; Riloff, E., Chiang, D., Hockenmaier, J., Tsujii, J., Eds.; pp. 4847–4853. [[CrossRef](#)]
30. Rudinger, R.; Naradowsky, J.; Leonard, B.; Van Durme, B. Gender Bias in Coreference Resolution. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), New Orleans, LA, USA, 1–6 June 2018; Walker, M., Ji, H., Stent, A., Eds.; pp. 8–14. [[CrossRef](#)]
31. Pratama, T.; Cahyadi, A. Effect of User Interface and User Experience on Application Sales. *IOP Conf. Ser. Mater. Sci. Eng.* **2020**, *879*, 012133. [[CrossRef](#)]
32. Rodríguez Cantelar, M.; Cal, D.d.l.; Estecha Garitagoitia, M.S. Grande Gutiérrez, A.; Martín, D.; Rodríguez Núñez-Milara, N.; Martínez Jiménez, R.; D'Haro Enríquez, L.F. Genuine2: An open domain chatbot based on generative models. In *Alexa Prize SocialBot Grand Challenge 4 Proceedings*; Amazon Science: New York, NY, USA, September 2021.
33. Adiwardana, D.; Luong, M.T.; So, D.R.; Hall, J.; Fiedel, N.; Thoppilan, R.; Yang, Z.; Kulshreshtha, A.; Nemade, G.; Lu, Y.; et al. Towards a Human-like Open-Domain Chatbot. *arXiv* **2020**, arXiv:2001.09977.
34. Roller, S.; Dinan, E.; Goyal, N.; Ju, D.; Williamson, M.; Liu, Y.; Xu, J.; Ott, M.; Shuster, K.; Smith, E.M.; et al. Recipes for building an open-domain chatbot. *arXiv* **2020**, arXiv:2004.13637.
35. Xu, J.; Szlam, A.; Weston, J. Beyond Goldfish Memory: Long-Term Open-Domain Conversation. *arXiv* **2021**, arXiv:2107.07567.
36. Komeili, M.; Shuster, K.; Weston, J. Internet-Augmented Dialogue Generation. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Dublin, Ireland, 22–27 May 2022; Muresan, S., Nakov, P., Villavicencio, A., Eds.; pp. 8460–8478. [[CrossRef](#)]
37. Lee, J.; Shim, M.; Son, S.; Park, C.; Kim, Y.; Lim, H. There is no rose without a thorn: Finding weaknesses on BlenderBot 2.0 in terms of Model, Data and User-Centric Approach. *arXiv* **2022**, arXiv:2201.03239.
38. Kobza, O.; Herel, D.; Cuhel, J.; Gargiani, T.; Pichl, J.; Marek, P.; Konrad, J.; Sedivy, J. Enhancements in BlenderBot 3: Expanding Beyond a Singular Model Governance and Boosting Generational Performance. *Future Internet* **2023**, *15*, 384. [[CrossRef](#)]
39. Khatri, C.; Hedayatnia, B.; Venkatesh, A.; Nunn, J.; Pan, Y.; Liu, Q.; Song, H.; Gottardi, A.; Kwatra, S.; Pancholi, S.; et al. Advancing the State of the Art in Open Domain Dialog Systems through the Alexa Prize. *arXiv* **2018**, arXiv:1812.10757.
40. Gopalakrishnan, K.; Hedayatnia, B.; Chen, Q.; Gottardi, A.; Kwatra, S.; Venkatesh, A.; Gabriel, R.; Hakkani-Tür, D. Topical-Chat: Towards Knowledge-Grounded Open-Domain Conversations. *Proc. Interspeech* **2019**, 1891–1895. [[CrossRef](#)]
41. Rashkin, H.; Smith, E.M.; Li, M.; Boureau, Y.L. Towards Empathetic Open-domain Conversation Models: A New Benchmark and Dataset. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Florence, Italy, 28 July–2 August 2019; Korhonen, A., Traum, D., Màrquez, L., Eds.; pp. 5370–5381. [[CrossRef](#)]
42. Radlinski, F.; Balog, K.; Byrne, B.; Krishnamoorthi, K. Coached Conversational Preference Elicitation: A Case Study in Understanding Movie Preferences. In Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue, Stockholm, Sweden, 11–13 September 2019; Nakamura, S., Gasic, M., Zuckerman, I., Skantze, G., Nakano, M., Papangelis, A., Ultes, S., Yoshino, K., Eds.; pp. 353–360. [[CrossRef](#)]
43. Hu, E.J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; Chen, W. LoRA: Low-Rank Adaptation of Large Language Models. In Proceedings of the International Conference on Learning Representations, Virtual Event, 25–29 April 2022.
44. Soltan, S.; Ananthakrishnan, S.; FitzGerald, J.G.M.; Gupta, R.; Hamza, W.; Khan, H.; Peris, C.; Rawls, S.; Rosenbaum, A.; Rumshisky, A.; et al. AlexaTM 20B: Few-shot learning using a large-scale multilingual seq2seq model. *arXiv* **2022**, arXiv:2208.01448.
45. Wolf, M.J.; Miller, K.W.; Grodzinsky, F.S. Why we should have seen that coming: Comments on microsoft's tay "experiment", and wider implications. *Orbit J.* **2017**, *1*, 1–12. [[CrossRef](#)]
46. Elsafoury, F.; Katsigiannis, S.; Pervez, Z.; Ramzan, N. When the Timeline Meets the Pipeline: A Survey on Automated Cyberbullying Detection. *IEEE Access* **2021**, *9*, 103541–103563. [[CrossRef](#)]
47. Sun, H.; Xu, G.; Deng, J.; Cheng, J.; Zheng, C.; Zhou, H.; Peng, N.; Zhu, X.; Huang, M. On the Safety of Conversational Models: Taxonomy, Dataset, and Benchmark. In Proceedings of the Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, 22–27 May 2022; pp. 3906–3923. [[CrossRef](#)]
48. Nadeem, M.; Bethke, A.; Reddy, S. StereoSet: Measuring stereotypical bias in pretrained language models. *arXiv* **2020**, arXiv:2004.09456.
49. Bojanowski, P.; Grave, E.; Joulin, A.; Mikolov, T. Enriching Word Vectors with Subword Information. *Trans. Assoc. Comput. Linguist.* **2017**, *5*, 135–146. [[CrossRef](#)]

- 
50. Le, C.P.; Dai, L.; Johnston, M.; Liu, Y.; Walker, M.; Ghanadan, R. Improving open-domain dialogue evaluation with a causal inference model. *arXiv* **2023**, arXiv:2301.13372.
  51. Antypas, D.; Ushio, A.; Camacho-Collados, J.; Silva, V.; Neves, L.; Barbieri, F. Twitter Topic Classification. In Proceedings of the 29th International Conference on Computational Linguistics, Gyeongju, Republic of Korea, 12–17 October 2022; pp. 3386–3400.

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.