

# SAPIEN: Affective Virtual Agents Powered by Large Language Models\*

Masum Hasan\*, Cengiz Ozel†, Sammy Potter‡ and Ehsan Hoque§

Department of Computer Science, University of Rochester

Rochester, NY, United States

Email: {\*m.hasan@, †cozel@cs., ‡spotter14@u., §mehoque@cs.} rochester.edu

**Abstract**—In this demo paper, we introduce SAPIEN, a platform for high-fidelity virtual agents driven by large language models that can hold open domain conversations with users in 13 different languages, and display emotions through facial expressions and voice. The platform allows users to customize their virtual agent’s personality, background, and conversation premise, thus providing a rich, immersive interaction experience. Furthermore, after the virtual meeting, the user can choose to get the conversation analyzed and receive actionable feedback on their communication skills. This paper illustrates an overview of this technology, ranging from entertainment to mental health, communication training, language learning, education, healthcare, and beyond. Additionally, we consider the ethical implications of such realistic virtual agent representations and the potential challenges in ensuring responsible use.

**Index Terms**—Virtual Avatars, Virtual Agents, Affective AI, Large Language Models

## I. INTRODUCTION

Allowing a user to define the traits and characteristics of a virtual agent, carrying a dynamic conversation, and receiving automated feedback has been an open-ended research problem for many years [1]. The rapid advancement of Large Language Models (LLMs) in recent years has enabled possibilities in designing user experiences that didn’t exist before [2]–[4]. In this demo, we present Synthetic Anthropomorphic Personal Interaction ENgine (SAPIEN), a platform for LLM-powered high-fidelity virtual agents that can engage in real-time open-domain conversations, while also expressing emotions through voice and facial expressions.

One of the notable features of SAPIEN is its extensive range of customization options, allowing users to engage in immersive and meaningful interactions. Users can choose from a wide range of virtual agent avatars that reflect a diverse array of ages, gender, and ethnicities. Going further, users can select the desired personality, background, and conversational context of a virtual agent, creating an experience tailored to their specific needs or preferences.

Once a virtual agent is selected and its traits are defined, users can begin a real-time video call interaction with it. With the help of the large language model, the virtual agents dynamically adjust their emotional state, vocal, and facial expressions, showcasing a spectrum of seven basic emotions.

NSF and NSF REU IIS-1750380, Seedling from Goergen Institute for Data Science (GIDS), and Gordon and Moore Foundation.



Fig. 1. Face-to-face video call interaction with SAPIEN™ Virtual Agent

SAPIEN leverages state-of-the-art models in Speech-to-Text [5], [6], Text-to-Speech [7]–[9], and large language modeling [2], [4], [10]–[14]. The virtual agents fluently speak thirteen different languages and counting, making it accessible across a global user base.

Upon finishing a video call with the virtual agents, a user can choose to get their conversation analyzed for personalized feedback. The system provides AI-generated feedback to the user based on the user’s goal. The user can decide the topic of the feedback to suit their learning goal and repeat the conversation until the learning goal is met. The inherent flexibility of the virtual agent persona and the feedback could make it potentially applicable to a myriad of applications, including communication training, language learning, and professional applications like healthcare, sales, and leadership training.

With the rising technical capabilities of LLMs, there is expected to be a drastic shift in the labor market in the coming years [15]. According to recent studies [15], the importance of the job market is going to shift from hard technical skills to soft “human” skills. In this changing landscape, SAPIEN can help people adapt and cope, by helping them cultivate human skills with the help of AI.

## II. SYSTEM DESCRIPTION

The overall working of SAPIEN Virtual Agents, referred to as ‘Bot’ for simplicity, is represented in Figure 2. The SAPIEN system is initialized when a user’s speech utterance is captured and transmitted to our back-end server for processing. This utterance is transcribed into text by a high-precision Speech

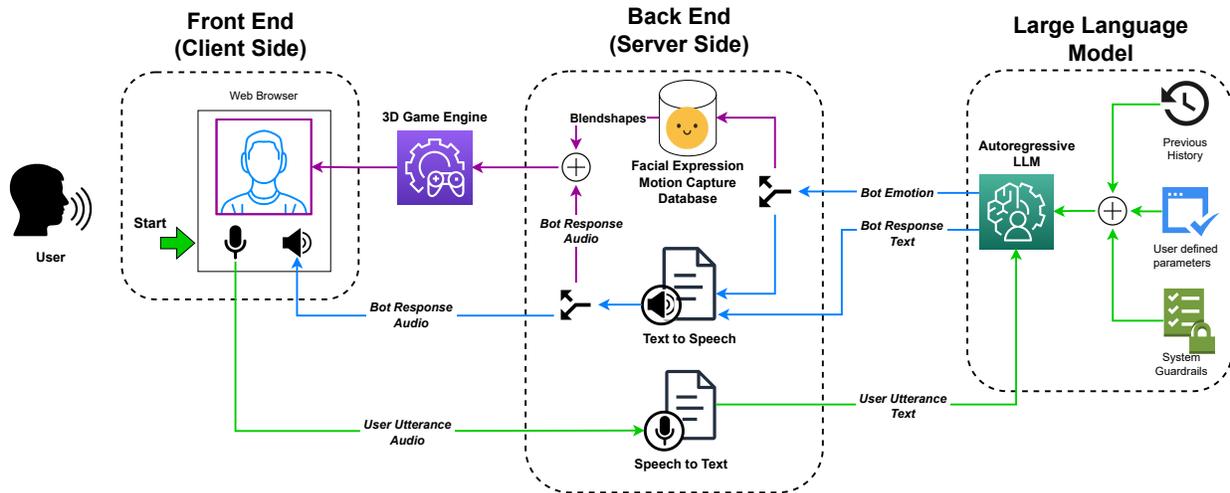


Fig. 2. A single turn conversation flow in SAPIEN. User utterance is transcribed and sent to LLM. The LLM response is spoken out by the virtual agent.

to Text (STT) model [5], [6], [16], [17] and subsequently processed by an autoregressive Large Language Model (LLM) fine-tuned for instruction following [3], [4], [10]–[14], [18].

The LLM is conditioned on user-defined parameters like personality traits, conversation premise, user information, and previous conversation history. To prevent inappropriate or offensive behavior, the LLM also adheres to system guardrails. A notable aspect of the LLM is also predicting the virtual agent’s emotional state. Conditioning on the user-defined parameters, system guardrails, and previous conversation history, the LLM is instructed to generate the bot’s response, alongside the appropriate emotional state of the bot from the following list: Neutral, Happy, Sad, Angry, Surprised, Afraid, and Disgusted.

This emotional state, along with the text response, is used to generate an audio file of the bot’s response using a Text to Speech (TTS) model. Concurrently, the emotional state triggers the selection of a corresponding facial expression from our pre-recorded motion capture database. This facial expression data, in the form of blendshapes, is passed to a 3D game engine to animate the virtual agent.

The resultant animation and generated audio are synchronized, forming a coherent, visually expressive response from the virtual agent. This combined output is streamed to the user’s web browser in near real-time, allowing for an immersive experience close to an actual video call.

Once the conversation is over, the user can opt-in to receive feedback on their conversation. An LLM is instructed to analyze the conversation transcript based on the user’s goal, identify strengths and weaknesses on the user’s communication skill, and generate actionable feedback for the user.

### III. APPLICATIONS

The customizability of the conversation scenario, dynamic dialogues, and the feedback system combined make SAPIEN uniquely suitable for a variety of communication training purposes. For example, the system can be used as a com-

munication practice tool for people with social anxiety or neurodiversity [19], [20], public speaking [21], job interviews [22], helping elderly with social skills [23], and even speed dating [24]. It also has an excellent potential for professional applications. Such as training doctors in bedside manners or delivering difficult news to their patients [25], personalized training for leadership, business negotiation, sales, marketing, etc. The multilingual ability makes the platform a powerful tool for language learners. Furthermore, the non-judgemental, low stake, repeatable conversations with virtual agents make the platform a helpful tool for anyone to roleplay any difficult interpersonal scenario in a personal or professional setup.

### IV. THE DEMO

Our platform is hosted in the cloud and accessible from any part of the world. During the conference demo, we wish to have the visitors live interact with SAPIEN virtual agents in a variety of interesting scenarios and receive immediate feedback on their communication skills. We will also prepare some pre-recorded user interaction videos to demonstrate any rare or difficult cases or as a backup for technical failures.

### ETHICAL IMPACT STATEMENT

SAPIEN is designed to augment and enrich our capacity for communication, empathy, and understanding, but not substitute human connections. To safeguard against potential emotional dependencies on the system, SAPIEN does not retain the memory of previous interactions, and the conversations are limited to a 10 minutes window with a warning at the 8-minute mark. To prevent the practice of bullying or abusive behaviors using our system, we enabled our virtual agents to end the video call if the user repeatedly displays aggressive or offensive behavior. We are continuously investigating more safety and ethical issues regarding the use of the system.

## REFERENCES

- [1] M. E. Hoque and R. W. Picard, "Rich nonverbal sensing technology for automated social skills training," *Computer*, vol. 47, no. 4, pp. 28–35, 2014.
- [2] OpenAI, "Introducing chatgpt," <https://openai.com/blog/chatgpt>, (Accessed on 06/22/2023).
- [3] "Anthropic — introducing claude," <https://www.anthropic.com/index/introducing-claude>, (Accessed on 06/22/2023).
- [4] G. AI, "An important next step on our ai journey," 2023. [Online]. Available: <https://blog.google/technology/ai/bard-google-ai-search-updates/>
- [5] J. Li, "Recent advances in end-to-end automatic speech recognition," *APSIPA Transactions on Signal and Information Processing*, April 2022. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/recent-advances-in-end-to-end-automatic-speech-recognition/>
- [6] W. Xiong, L. Wu, F. Alleva, J. Droppo, X. Huang, and A. Stolcke, "The microsoft 2017 conversational speech recognition system," in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2018, pp. 5934–5938.
- [7] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio *et al.*, "Tacotron: Towards end-to-end speech synthesis," *arXiv preprint arXiv:1703.10135*, 2017.
- [8] R. Luo, X. Tan, R. Wang, T. Qin, J. Li, S. Zhao, E. Chen, and T.-Y. Liu, "Lightspeech: Lightweight and fast text to speech with neural architecture search," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 5699–5703.
- [9] S.-g. Lee, H. Kim, C. Shin, X. Tan, C. Liu, Q. Meng, T. Qin, W. Chen, S. Yoon, and T.-Y. Liu, "Priorgrad: Improving conditional denoising diffusion models with data-driven adaptive prior," *ICLR*, 2022.
- [10] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto, "Stanford alpaca: An instruction-following llama model," [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- [11] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [12] OpenAI, "Gpt-4 technical report," 2023.
- [13] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, "Training language models to follow instructions with human feedback," *Advances in Neural Information Processing Systems*, vol. 35, pp. 27 730–27 744, 2022.
- [14] A. Köpf, Y. Kilcher, D. von Rütte, S. Anagnostidis, Z.-R. Tam, K. Stevens, A. Barhoum, N. M. Duc, O. Stanley, R. Nagyfi *et al.*, "Openassistant conversations—democratizing large language model alignment," *arXiv preprint arXiv:2304.07327*, 2023.
- [15] T. Eloundou, S. Manning, P. Mishkin, and D. Rock, "Gpts are gpts: An early look at the labor market impact potential of large language models," *arXiv preprint arXiv:2303.10130*, 2023.
- [16] Y. Leng, X. Tan, L. Zhu, J. Xu, R. Luo, L. Liu, T. Qin, X. Li, E. Lin, and T.-Y. Liu, "Fastcorrect: Fast error correction with edit alignment for automatic speech recognition," *Advances in Neural Information Processing Systems*, vol. 34, pp. 21 708–21 719, 2021.
- [17] W. Hou, J. Wang, X. Tan, T. Qin, and T. Shinozaki, "Cross-domain speech recognition with unsupervised character-level distribution matching," *INTERSPEECH*, 2021.
- [18] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez, I. Stoica, and E. P. Xing, "Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality," March 2023. [Online]. Available: <https://lmsys.org/blog/2023-03-30-vicuna/>
- [19] M. R. Ali, S. Z. Razavi, R. Langevin, A. Al Mamun, B. Kane, R. Rawassizadeh, L. K. Schubert, and E. Hoque, "A virtual conversational agent for teens with autism spectrum disorder: Experimental results and design lessons," in *Proceedings of the 20th ACM International Conference on Intelligent Virtual Agents*, ser. IVA '20. New York, NY, USA: Association for Computing Machinery, 2020. [Online]. Available: <https://doi.org/10.1145/3383652.3423900>
- [20] S. Z. Razavi, M. R. Ali, T. H. Smith, L. K. Schubert, and M. E. Hoque, "The lissa virtual human and asd teens: An overview of initial experiments," in *Intelligent Virtual Agents*, D. Traum, W. Swartout, P. Khooshabeh, S. Kopp, S. Scherer, and A. Leuski, Eds. Cham: Springer International Publishing, 2016, pp. 460–463.
- [21] M. Fung, Y. Jin, R. Zhao, and M. E. Hoque, "Roc speak: Semi-automated personalized feedback on nonverbal behavior from recorded videos," in *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, ser. UbiComp '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 1167–1178. [Online]. Available: <https://doi.org/10.1145/2750858.2804265>
- [22] M. E. Hoque, M. Courgeon, J.-C. Martin, B. Mutlu, and R. W. Picard, "Mach: My automated conversation coach," in *Proceedings of the 2013 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, ser. UbiComp '13. New York, NY, USA: Association for Computing Machinery, 2013, p. 697–706. [Online]. Available: <https://doi.org/10.1145/2493432.2493502>
- [23] S. Z. Razavi, L. K. Schubert, K. van Orden, M. R. Ali, B. Kane, and E. Hoque, "Discourse behavior of older adults interacting with a dialogue agent competent in multiple topics," *ACM Trans. Interact. Intell. Syst.*, vol. 12, no. 2, jul 2022. [Online]. Available: <https://doi.org/10.1145/3484510>
- [24] M. R. Ali, D. Crasta, L. Jin, A. Baretto, J. Pachter, R. D. Rogge, and M. E. Hoque, "Lissa — live interactive social skill assistance," in *2015 International Conference on Affective Computing and Intelligent Interaction (ACII)*, 2015, pp. 173–179.
- [25] M. R. Ali, T. Sen, B. Kane, S. Bose, T. M. Carroll, R. Epstein, L. Schubert, and E. Hoque, "Novel computational linguistic measures, dialogue system and the development of sophie: Standardized online patient for healthcare interaction education," *IEEE Trans. Affect. Comput.*, vol. 14, no. 1, p. 223–235, jan 2023. [Online]. Available: <https://doi.org/10.1109/TAFFC.2021.3054717>