# 2023 年度ウェブプログラミング実習 総合実習レポート

# Hopper Hug

### 1. 作成概要

スマートフォンを使用したインタラクションなくまのぬいぐるみ(図 1)を制作した。ユ ーザを記憶して友達同士のようなコミュニケーションを楽しむことができるぬいぐるみで ある。ユーザ登録をすれば誰でもゆったりとした会話が楽しむことができる。小さい子供だ けではなく、あらゆる世代の方に使ってほしいシステムである。

本システムでは入力された発話に対してぬいぐるみが返答をすることが可能である。ユ ーザごとの過去の対話情報も記憶しているため、機械的ではない人間らしい会話をするこ とが可能である。またカメラ機能も搭載しており、カメラからの視覚的情報を基に会話をす ることも可能である。他にも人間のハグを検知して話しかけてくれるハグ機能や、ぬいぐる みに危害が加わった際に助けを求めるヘルプ機能も搭載している。

システムの概要を図1に示す。データベースには、ユーザのユーザ名、パスワード、ロ グイン履歴、過去の対話情報などが保存されており、PHPからデータベースに接続し、ユ ーザー登録処理、ログイン処理、会話履歴の記録と検索を行っている。実行画面では、 JavaScriotを用いて、音声認識、カメラのオンオフの切り替え、ヘルプ機能に使用するス マートフォンの加速度センサの値の取得を行っている。

この他に obniz と WebAPI を用いて会話機能やハグ機能を実現した。obniz の人感セン サを用いてハグ機能の判定を行った。また、ChatGPT で入力された発話に対する返答を生 成し、その返答を VOICEVOX で音声データに変換して再生した。



図 1 Hopper Hug の外観



図 2 システム概要図

2.	担当箇所について
----	----------

表 1 担当表

項目	担当者	備考
データベース	А	データベースと SQL 文の設計
サーバーサイド	А	PHP を用いたユーザ管理、会話履歴管理
会話生成	B、本間	WebAPI (GPT3, GPT4-V)を用いた会話の生成
会話の要約	В	WebAPI (GPT3)を用いて記録用に会話を要約
音声認識、合成	本間	APIを用いて音声を検出&認識し、返答音声を合成
カメラ機能	本間	Web カメラの情報をもとにして会話を生成
加速度センサ	С	JavsScript でスマホの加速度センサの値を取得
IOT(Obniz)	С	JavsScript で Obniz と接続し、センサの値を取得
デザイン	A、本間	CSS + JavaScript によるアニメーション
プレゼン	全員	パワーポイントの作成、発表、動画の作成

担当表を表1に示す。著者の担当は、音声認識と合成、カメラ機能、会話生成、デザインである。また、それぞれメンバーが作成してきたコードを最終的に合体する作業も著者が担当した。本章では、ユーザーとHopperHugが会話する流れを、I音声認識、II webカメラ撮影、III会話生成、IV音声合成に分けて、それぞれ実装を説明していく。最後 に、キャラクターの会話状況を示すアニメーションについて説明する。



図 3 index.php の初期画面

### I. 音声認識

ユーザーは、ログイン後、index.php(図3)で会話を行う。Startを押すと音声 認識が始まる。音声認識には、ウェブ音声認識 API である Webkit Speech Recognitionを用いている(参考文献3)。この API を用いることで、ブラウザが対 応していればデバイスに備え付けられた音声認識をシンプルなコードで利用するこ とができる。Webkit Speech Recognitionのメンバ変数、onresultにコールバック 関数を代入することで、認識結果を取得している。ユーザーが喋るのを止めたこと を検知すると、onspeechend に登録されたコールバック関数が呼ばれる。このコー ルバック関数内で、Ⅱのカメラ撮影やⅢの返答合成を行っている。

また、IVの音声合成と再生が終了次第、Webkit Speech Recognition オブジェクトの start()関数を呼び、音声認識を再開させている。これにより、ユーザーは連続した会話が可能である。また、再生された合成音声を認識してしまうこともこれにより防いでいる。なお、音声認識関連のコードは src/script/main. js に記述している。

II. Web カメラ撮影

図3左上のトグルスイッチを押すことで、カメラ映像の入力(以降、おめめモード)の 0N/0FF を切り替えることができる。このトグルスイッチは参考文献2を参考に作成した。ここで、カメラ自体の起動は、スタートボタンを押した後に行っている。これは、Safari などのセキュリティに厳しい一部のブラウザでは、カメラの起動がユーザーのジェスチャ(ボタンを押すなどの挙動)によって行われなければいけないためである。Start ボタンが押されると、カメラが起動し、

style="display:none"属性が適応された video タグオブジェクトにインカメラの 映像を映し出す。インカメラの映像は、

video.srcObject =

await navigator. mediaDevices. getUserMedia({ video: { facingMode: "user" } }) のように、facingMode に"user"を指定することで取得している。 I にて、音声認 識が終わり、おめめモードが ON の場合は、video タグに映し出されている映像を canvas に drawImage 関数を用いて書き込み、canvas. toDataURL 関数を用いて jpeg の base64 エンコードを行い、URL を取得する。この URL は、IIIで使用される。ま た、コードは src/script/camera. js に記述している。

#### III. 会話生成

I で認識したテキストと、II で取得された画像 URL をもとに適切な返答を OpenAI の API を用いて生成する。モデルは、おめめモードが ON の場合は、画像入力に対応 している GPT4-Vision を、OFF の場合は、GPT-3.5-turbo を用いている。GPT4-Vision は返答に時間がかかる傾向にあるため、このように使い分けている。fetch 関数を用いてリクエストを行っている。リクエストのパラメータは OpenAI の公式サ イト(参考文献 4)を参考にした。また、以下のようにプロンプトを設定した。

"あなたはクマのキャラクターで名前は花子です、ユーザーと会話をしてください。ユーザーの名前は、 \${user\_name}です. 会話はかわいいキャラクターのようにしてください。後、一人称を使用する際は一人称 を「ぼく」にして、返答は短めにしてください。会話が続かなさそうであったら自ら話題を振りなさい. ま た,以下はあなたと\${user\_name}の今までの会話内容を要約したものです. 適宜この内容を参照しながら会 話をしてください. 無理に使う必要はありません 返答は短めにすることをこころがけてください. 返答は短 くまとめ,レスポンスには返答の内容のみ入れてください.

<<以下は、お目々モード ON のとき、挿入>>

必要であれば、適度に視覚的情報を用いて返答をしなさい、無理に視覚情報を用いる必要はありません.なお 与えられた写真はあなたの目に入っている情報とします。写真とは言わないでください"

例:

ユーザー:いい天気だー あなたの返答例: いい天気だね!花子といっしょに遊ぼう!\${user\_name}さん! \${record list}

User\_name には、PHP によって html に書き出されている、ログイン中のユーザー 名が取得され、代入される。また、record\_list には、そのユーザーの会話履歴が 文字列化され代入される。このようなプロンプトを用いることで、可愛いくまのキ ャラクターを演じさせ、会話履歴を踏まえた会話をさせている。また、おめめモー ドが ON の場合は、web カメラの情報も含めて答えさせることで、視覚的情報を交え た会話が可能となる。ここで生成された返答は、IVで利用される。実装のコード は、src/script/GPT. js にある。

#### IV. 音声合成、再生

Ⅲで生成された返答は、有志の方が公開している、VoiceVox を利用して音声を合 成できる API「TTS QUEST V3 VOICEVOX API」(参考文献 1)を利用して、音声 に変換している。参考文献 1 のコードを 1 部改変して、API にリクエストを送り、 合成が完了し次第、再生している。再生終了後には、音声認識が開始され、再び I に戻る。また、obniz が人感センサの反応を検知した場合には、「うれしい」という 音声を再生し(ハグ機能)、スマホの加速度センサが反応した場合には、「助けて」 という音声を再生する(ヘルプ機能)

ここで、Audio オブジェクトを新規作成するのではなく、既存の Audio オブジェ クトのソース(src 属性)を張り替えるという手法で再生を行っている。これは、モダ ンな Web ブラウザではセキュリティ対策によって、任意のタイミングで play()関数 を用いて音源を再生できないためである。II のカメラ利用のように、モダンブラウ ザでは、ユーザのジェスチャに由来した音楽の再生でないと再生がブロックされて しまう。これを避けるために、Start ボタンが押された段階で、autoplay 属性を true にした、Audio オブジェクトを作成し、空のデータをソース(src 属性)に入れて 再生しておく。そして、再生するタイミングで適切なソースに張り替えることでユ ーザーの介入なしに任意のタイミングでの再生を可能にしている。

音声を合成するコードは、src/script/tts.js に、音声を再生するコードは、 src/script/audioManager.js にある。

以上が、会話における実装の流れである。次に、実行状態を表すアニメーションの実装 について説明する。アニメーションの流れを図4に示す。Startを押すと音声認識中に なり、円が緑になり拡大、縮小をゆっくりと繰り返す。認識が終わり、会話が生成され ている途中では、色がオレンジになり、円が縮小して縦にジャンプをする。合成音声の 再生中は、色がピンクに変わりその場で回転をする。ユーザーは、今、どの処理が行わ れているのかを視覚的に知ることができる。アニメーションの実装には anim. js(参考文 献 5)を利用した。コードは src/script/animation.js にある。最後に、会話が終わったら 会話終了ボタンを押すことで、今までの会話が GPT によって要約されて、textarea に 自動で入力される。その後、自動で form が submit され、record.php によってデータ ベースに会話履歴が書き込まれる。



図 4 アニメーションの流れ

## 参考文献

1) ts-klassen, TTS QUEST V3 VOICEVOX API, (https://github.com/ts-klassen/ttsQuestV3Voicevox, (2024年1月参照)

2) W&M de Asobo, 【超簡単!】CSS と JavaScript でトグルスイッチを作ろう!, (https://web-de-asobo.net/2023/04/24/toggle-switch/, (2024年1月参照)

3) mozilla, ウェブ音声 API, https://developer.mozilla.org/ja/docs/Web/API/Web\_Speech \_API, (2024 年 1 月参照)

 $\label{eq:openAI} \ensuremath{\text{API}} \ensuremath{\text{Reference}} - \ensuremath{\text{Chat}}, \ensuremath{\,\text{https://platform.openai.com/docs/api-reference/chat} \\ \ensuremath{\text{Chat}} \ensuremath{\text{API}} \ensuremath{\text{Reference}} - \ensuremath{\text{Chat}} \ensuremath{\text{Chat}} \ensuremath{\text{Chat}} \ensuremath{\text{API}} \ensuremath{\text{Reference}} \ensuremath{\text{Chat}} \ensuremath{\text{Chat}} \ensuremath{\text{API}} \ensuremath{\text{Reference}} \ensuremath{\text{Chat}} \ensuremath{\text{Reference}} \ensuremath{\text{Chat}} \ensuremath{\text{API}} \ensuremath{\text{Reference}} \ensuremath{\text{Reference}} \ensuremath{\text{Chat}} \ensuremath{\text{Reference}} \en$ 

t, (2024年1月参照)

5) juliangarnier, Anime, js, <u>https://animejs.com/</u>