

# 対話型鑑賞のファシリテータ養成のための複合現実 AI システム

○福田 涼太 (九州大学), 倉爪 亮 (九州大学)

## Mixed reality AI system for training VTS facilitators

○ Ryota FUKUDA (Kyushu University), and Ryo KURADUME (Kyushu University)

Abstract: Visual Thinking Strategies, or VTS for short, is a method of art appreciation in which multiple people in a group interpret a work of art through repeated dialogue. In this paper, we propose a mixed reality AI system that uses MR headsets and Large Language Models to train facilitators, who play an important role in VTS.

### 1. 序論

美術鑑賞の方法の一つとして、複数人のグループで対話を重ねながら美術作品を読み解いていく対話型鑑賞 [1][2] という鑑賞法がある。対話型鑑賞のイメージ図を Fig.1 に示す。この対話型鑑賞では、グループ内の会話をスムーズに進行させるファシリテータの役割が重要となる。しかし、ファシリテータの訓練、養成には、現状では美術館等で開催される絵画鑑賞会などで、実際の参加者を対象に現場で練習を繰り返すしかなく、訓練のための十分な機会が確保できないという課題があった。そこで、我々は拡張現実ヘッドセットと大規模言語モデルを用いて対話型鑑賞を仮想的に再現することで、美術館等の場所の制約なく、一人で対話型鑑賞のファシリテータ訓練を行うことができる複合現実 AI システムを開発した。



Fig. 1 対話型鑑賞のイメージ (ChatGPT-4o にて作成)

### 2. システム構成

#### 2.1 対話型鑑賞についての状況設定

まず仮想的に再現する対話型鑑賞の状況を Fig.2 で示す。複合空間内には仮想的に5人のキャラクターが存在しており、それらは全員対話型鑑賞の参加者である。拡張現実ヘッドセットを装着している人はファシリテータとして振る舞い、Fig.2 の上部中央 (ユーザの正面)、および Fig.3 で示すように、ユーザの背面 (実際の対話型鑑賞と同じく鑑賞者の正面) にある絵画に対して対話型鑑賞を展開していくことが要求される。



Fig. 2 対話型鑑賞の状況

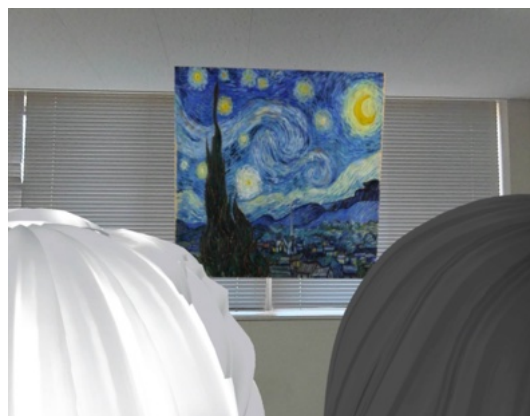


Fig. 3 ユーザの背面すなわち鑑賞者の正面にある絵画

#### 2.2 5人のキャラクターについて

各キャラクターはそれぞれ個別に感情の値を0から1の間の値で持っており、会話がポジティブな内容であれば値が大きくなり、ネガティブな内容であれば値が小さくなるというように、ユーザとの会話の内容に応じて値は徐々に変化していく。また、この感情の値により顔の表情も変化する。ヘッドセット (Quest3、Meta) のアイトラッキング機能を利用することでキャラクターと目が合っているかどうかを判定し、目が合っているキャラクターの表情にその時点での感情の値を反映させている。さらに、ファシリテータが各参加者に均等に話しかけるようにするために、キャラクターごとに目が合っている時間を計測して表示し、最も目が合っている時間が短いキャラクターには”short”と表示している (Fig.4)。Fig.2 の下部に表示されているように5人のキャラクターにはそれぞれ名前があり、年齢や出身、性格なども個別に設定されている。この

ような性格などの違いにより各キャラクターの発言内容や話し方が異なるようにしているが、詳細は 2.3 で説明する。



Fig. 4 目が合っている時間の計測

### 2.3 複合現実内のキャラクターとの会話の実現について

ユーザとキャラクターの会話をどのように実現しているかについて説明する。ユーザが発言してからキャラクターが返答するまでの大まかな流れを Fig.5 に示す。

まず対話型鑑賞において、目を合わせながら会話をすることが大切である。したがって目を見て発言しなければキャラクターが反応しない仕様にしており、目を見て発言することが会話を成立させる条件の一つ目となる。しかし、実際の対話型鑑賞ではずっと目を見て話すわけではなく、絵画に視線を向けたり、他の人の様子を見たりする。このような状況を再現するために、発言中はユーザは視線を動かすことができるが、発言が終わるタイミングで返答して欲しいキャラクターの目を見ている必要がある。会話を成立させる条件の二つ目は、一定の大きさ以上の音量で発言することである。わずかな音量であれば反応せず、一定の大きさ以上の音量であれば録音を開始し、無音状態が閾値以上続いた時点で録音を終了する。ここで、録音の長さが閾値以下であった場合はキャラクターは反応しない。つまり一定時間より長く発言するということが会話を成立させる条件の三つ目である。これら三つの条件を満たした場合に次のステップに移行する。

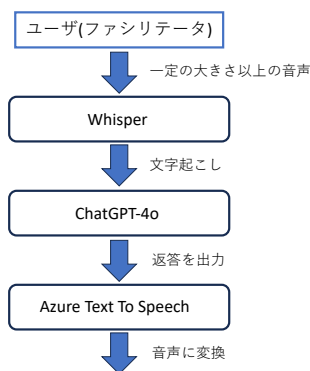


Fig. 5 ユーザの発言～キャラクターの返答までの流れ

#### 2.3.1 音声認識モデル Whisper

音声認識モデル Whisper(OpenAI) を利用し、ユーザが発言した音声を文字起こしして、テキストデータに変換する。

#### 2.3.2 大規模言語モデル GPT-4o

前のステップで得られたテキストデータを大規模言語モデル GPT-4o に読み込ませる。この時、GPT-4o には、対象キャラクターなどに関する設定が記述されたプロンプト (Fig.6、Fig.7)、対象とする絵画、および実際に複数人に聞き取りを行って得られた対象の絵画に関する印象などをまとめたデータ (Fig.8) を読み込ませている。なお、前述のように、プロンプトでは、キャラクターごとに異なる性格などの特徴を設定している。このプロンプトを読み込ませることで、そのキャラクターに合わせた返答を得ることができる。さらに、Fig.8 で示した印象データをポジティブな内容のものとネガティブな内容のものに分け、ポジティブな性格のキャラクターにはポジティブな印象データを反映させ、ネガティブな性格のキャラクターにはネガティブな印象データを反映させることで、よりそのキャラクターの性格に合った返答を得られるようにしている。各キャラクターにつき、一つの絵画に対して読み込ませる印象データは 100 個前後である。また、対話型鑑賞では 1 対 1 のやり取りではなく、複数人での会話であるため、誰がどのようなことを言ったのかを各キャラクターが把握する必要がある。そこで、今までに誰がどのような発言をしたのかを記録し、記録された会話の流れも GPT-4o に読み込ませるようにしている。

#### 2.3.3 音声合成システム Azure Text To Speech

GPT-4o により出力されたテキストデータを音声機能を提供する Azure Text To Speech(Microsoft) を利用することによって、音声に変換する。この時合成される音声は年齢、性別、国籍で異なるようにしている。このような流れで音声出力されることによって、ユーザとキャラクターとの会話を実現される。

#### 2.3.4 キャラクターの口の動き

キャラクターが発言していることを視覚的にわかるようにするために、Oculus が提供するリップシンクライブラリである OVR Lip Sync を用いることによって、音声に対応して口を動かすようにしている。

```
### SETTING OF A CHATBOT (YOU) ###
Please refer to the following for character settings.
Character Settings
Name: Satomi Hamaguchi
Number(番号): 1
Gender: Female
Background: Artist
Age: 56 years old
Personality: Cheerful, active, kind and sensitive
Favorite food: Orange
Hobbies: Taking a picture
Hometown: Hokkaido
Habit of saying: ファンファン

Please refer to the following for setting up a conversation situation.
Conversation situation
Place: Dimly lit museum
Situation: Appreciating a painting in the hallway of the museum
Content: Communicating impressions of the painting
People: You(浜口サトミ), User, Character number 2(名前: 真田マサヒロ), Character number 3(名前: 松林カイ), Character number 4(名前: 大空ツバサ), Character number 5(名前: 安森ウメ)
```

Fig. 6 キャラクター設定のプロンプトの一部: 例①





ファシリテータがテクニックをうまく使いながら対話型鑑賞を進行するための訓練につながると期待される。



Fig. 10 表示される指示 (3 種類)

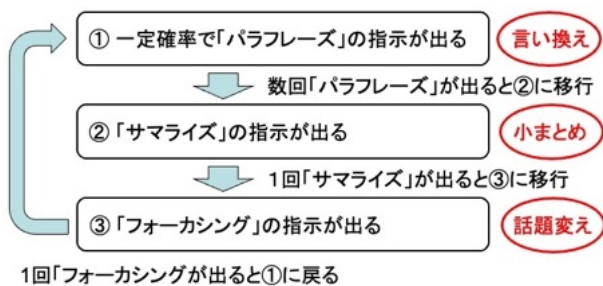


Fig. 11 表示される指示の流れ

### 2.4.8 Photo

この機能は前述の 6 種類の絵画・写真ではなく、現実空間に存在する絵画や風景などの写真を撮って取り込むことで、その写真に対して対話型鑑賞を行うものである。現段階ではスマートフォン等で撮った写真を専用の GitHub ページにアップロードし、アップロードした写真の URL を用いてシステムに取り込んでいる。今後は、ヘッドセット (Quest3) で直接写真を撮って取り込むなど、より使いやすい手法を検討している。

### 3. 今後の展開

今後このシステムの実践、評価を行い、その上で改良を行なっていく。現在は GPT-4o を利用しているが、Anthropic によって開発された AI チャットモデルである Claude と GPT-4o を比較し、どちらがこのシステムに適しているか検討したいと考えている。

### 4. 結論

対話型鑑賞においてファシリテータは重要な役割を担うが、ファシリテータの養成、訓練を行う機会が少ないという問題がある。そこで本稿では、現在我々が開発している対話型鑑賞のファシリテータ養成のための複合現実 AI システムを紹介した。本システムは、場所を問わず、一人でも対話型鑑賞におけるファシリテータの訓練を行うことができ、ファシリテータ養成、訓練の機会を増やすことができる。今後は、システムの評価実験を通して得られる評価、意見をもとにシステムを改良し、実際の対話型鑑賞にできるだけ近づける予定である。

### 謝辞

本研究は、JST 共創の場形成支援プログラム「共生社会をつくるアートコミュニケーション共創拠点」(JPMJPF2105) の支援を受けたものです。また印象データをご提供いただきました川畑秀明慶應義塾大学教授、開発にご協力いただきましたファシリテータの近藤乃梨子氏、春日美由紀氏に感謝いたします。

### References

- [1] Philip Yenawine, 京都造形芸術大学アートコミュニケーション研究センター：“学力をのばす美術鑑賞 ヴィジュアル・シンキング・ストラテジー：どこからそう思う?”, 淡交社, -, 2015
- [2] 稲庭彩和子, 伊藤達矢, 河野佑美, 鈴木智香子, 渡邊祐子：“こどもと大人のためのミュージアム思考”, 左右社, -, 2022