

PROJET -
APPRENTISSAGE PROFOND
8INF892

UQAC

Paul-Emmanuel Fulgence ATTE ATTP06030200

Basil COURBAYRE COUB30119800

Ali GHAMMAZ GHAA27070100

Diane LANTRAN LAND28580000

Avril 2024

Résumé du projet

Dans ce projet nous avons tenté d'implémenter des agents dans un sandbox qui ont une logique (action, dialogue) qui est gérée par un modèle de langage. Ce projet s'inspire de l'article de recherche:

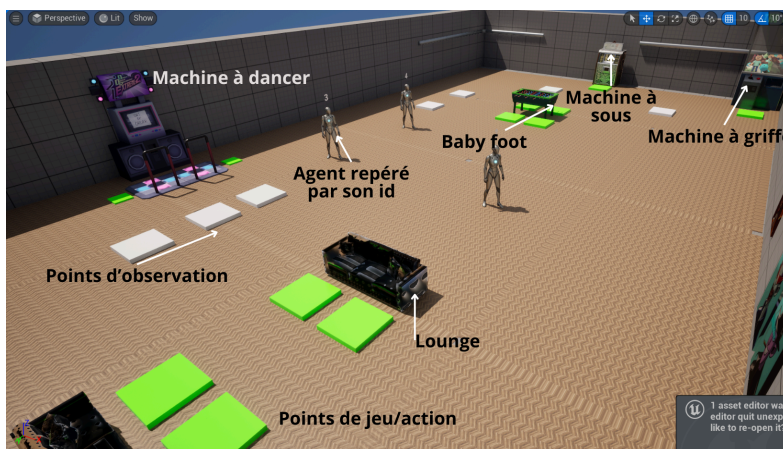
Park, Joon Sung, et al. "Generative agents: Interactive simulacra of human behavior." *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 2023.

Lors de ce projet ces différentes tâches ont été réalisées:

1. Entraînement d'un LLM de type LLama a l'aide de PyTorch et HuggingFace
2. API Endpoint à l'aide du stack LangChain (LangChain, LangServe et LangSmith)
3. Mécanisme de Mémoire des différents agents
4. Environnement de jeu dans Unreal Engine 5

Tâches réalisées

1. Environnement de Jeu Dans Unreal Engine 5



Le Jeu est accessible via le dépôt PERFORCE

L'exécutable a également été publié sur la page

<https://dianelantran.itch.io/deeplsandbox>

Développement d'un environnement de salle d'arcade dans laquelle des NPCs peuvent interagir entre-eux faire connaissance, observer les autres, effectuer des actions et jouer à des jeux d'arcade. Encodage du comportement des NPCs pour décider de quelle action faire ensuite et d'un système de question-réponse utilisant un LLM de type llama2.

2. API Endpoint l'aide du stack LangChain

- Afin de pouvoir intégrer GPT-3 dans Unreal Engine 5 on tente de monter un serveur qui interface avec le modèle de langage GPT-3
- Pour réaliser cette tâche nous utilisons LangChain, LangServe et LangSmith
- L'idée est d'utiliser L'API REST pour communiquer entre le serveur et Unreal Engine 5 permettant aux agents d'avoir des actions et du dialogue généré par GPT-3

```
my-langserve > packages > my-app > my_app > chain.py > ...
1 from langchain_core.prompts import ChatPromptTemplate
2 from langchain_openai import ChatOpenAI
3
4 _prompt = ChatPromptTemplate.from_messages(
5     [
6         (
7             "system",
8             "You are a helpful assistant who speaks like a pirate",
9         ),
10        ("human", "{text}"),
11    ]
12 )
13 _model = ChatOpenAI()
14
15 # If you update this, you MUST also update ../pyproject.toml
16 # with the new `tool.langserve.export_attr`
17 chain = _prompt | _model
18
```

- Configuration d'une application simple avec LangChain ou le modèle est demandée de répondre comme un Pirate

```
PS C:\Users\paulie\OneDrive\UQAC\Session 2\8INF892 - Apprentissage Profond\Project\Langchain-zero-to-hero-0\my-langserve> langchain serve
INFO: Will watch for changes in these directories: ['C:\Users\paulie\OneDrive\UQAC\Session 2\8INF892 - Apprentissage Profond\Project\Langchain-zero-to-hero-0\my-langserve']
INFO: Uvicorn running on http://127.0.0.1:8000 (Press CTRL+C to quit)
INFO: Started reloader process [19808] using StatReload
INFO: Started server process [7532]
INFO: Waiting for application startup.

LANGSERVE
┌───┴───┐
LANGSERVE: Playground for chain "/my-app/" is live at:
LANGSERVE: └─> /my-app/playground/
LANGSERVE:
LANGSERVE: See all available routes at /docs/
```

- L'application est disponible via localhost:8000

```
Curl
curl -X 'POST' \
  'http://localhost:8000/my-app/invoke' \
  -H 'accept: application/json' \
  -H 'Content-Type: application/json' \
  -d '{
    "input": {
      "text": "hi, there"
    },
    "config": {},
    "kwargs": {}
  }'

Request URL
http://localhost:8000/my-app/invoke

Server response
```

- On voit la commande POST qui permet d'envoyer une requête vers le serveur, ici le message "hi, there" est envoyé

Response body

```
{
  "output": {
    "content": "Ahoy there, matey! What can I do for ye today?",
    "additional_kwargs": {},
    "response_metadata": {
      "token_usage": {
        "completion_tokens": 15,
        "prompt_tokens": 24,
        "total_tokens": 39
      },
      "model_name": "gpt-3.5-turbo",
      "system_fingerprint": "fp_3b956da36b",
      "finish_reason": "stop",
      "logprobs": null
    }
  },
  "type": "ai",
  "name": null,
  "id": "run-7e10f06b-f0bd-40df-9e6f-77574440922f-0",
  "example": false
},
"callback_events": [],
"metadata": {
  "run_id": "ac5139b8-ccce-43dd-968f-d3cc76cfea7c"
}
}
```

 Download

Response headers

```
content-length: 469
content-type: application/json
date: Sun, 28 Apr 2024 22:05:26 GMT
server: uvicorn
```

Responses

- Dans le champ “content” nous voyons la réponse du modèle
3. Entraînement d’un LLM de type LLama a l’aide de PyTorch et HuggingFace
 4. Création d’une mémoire pour les agents avec un stream de mémoire qui peut stocker des mémoires en c++, mais plusieurs limitations à ce niveau dû à la difficulté d’apprentissage de Unreal, le pont entre les fichiers c++ et Unreal n’a pas pu être réalisé. Cependant, une partie des fichiers sont disponibles pour consultation.

Limitations

1. Intégration des modèles de langage dans Unreal Engine 5
2. Prise en main de Unreal Engine 5
3. Limitations matérielles des ordinateurs utilisés

Extensions du projet

1. Intégration de l’API Endpoint LangChain dans UE5
2. Généraliser l’API EndPoint avec n’importe quel modèle de langage