*User Question* → Select random component parts for chain → Generate text response → Present answer to user → User rates/scores answer → (back to) Select random component parts for chain
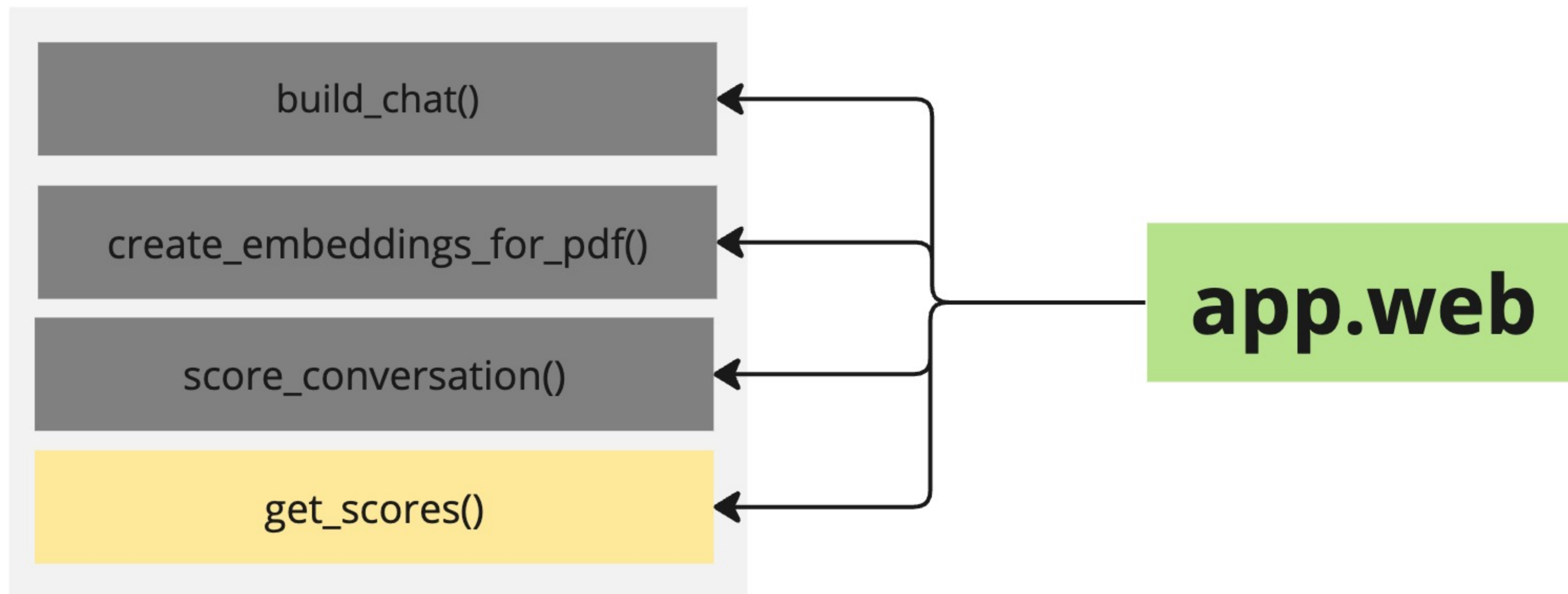
**app.chat**
**Module**

build_chat()

create_embeddings_for_pdf()

score_conversation()

get_scores()

**app.web**

# Browser

**+1** to conversation_id 123

**+0** to conversation_id 123

Server

## SQLite Database List of Conversations

| id | llm | memory | retriever |
|-----|-------|--------|-----------|
| 123 | gpt-4 | window | pinecone_2 |
| | | | |
| | | | |

# Redis

| Avg LLM Score | | Avg Retriever Score | | Avg Memory Score | |
|---------------|------|---------------------|-----|------------------|-----|
| gpt-4 | 0.75 | pinecone_2 | .4 | buffer | .8 |
| gpt-3.5-turbo | .2 | pinecone_3 | .1 | window | .4 |

*The way in which we store this data will be a little more complicated*

**Browser**

`< >`

+1 to conv '123'

+1 to conv '123'

+1 to conv '123'

+0 to conv '123'

+0 to conv '123'

Server

**SQLite Database List of Conversations**

| id | llm | memory | retriever |
|-----|-------|--------|------------|
| 123 | gpt-4 | window | pinecone_2 |
| | | | |
| | | | |

**Redis**

| **LLM Score Total** | | **Retriever Score Total** | | **Memory Score Total** | |
|---------------------|---|----------------------------|---|-------------------------|---|
| gpt-4 | 1 | pinecone_2 | | buffer | |
| gpt-2 | | pinecone_3 | | window | |

| **LLM Score Count** | | **Retriever Score Count** | | **Memory Score Count** | |
|---------------------|---|----------------------------|---|-------------------------|---|
| gpt-4 | | pinecone_2 | | buffer | |
| gpt-3.5-turbo | | pinecone_3 | | window | |

# Redis

### LLM Score Total

| | |
|---|---|
| gpt-4 | 1 |
| gpt-3.5-turbo | 5 |
| gpt-2 | 1 |

### LLM Score Count

| | |
|---|---|
| gpt-4 | 10 |
| gpt-3.5-turbo | 5 |
| gpt-2 | 3 |

## LLM Map

| | |
|---|---|
| gpt-4 | *builder* |
| gpt-3.5-turbo | *builder* |

Get all the values and counts from Redis

Take all the keys from the LLM map (avoid calculating average for old/removed LLM's)

For each active LLM, calculate average score. Use min of .1 to make sure a LLM always has a chance of being used

Calculate sum total of all ratings

Pick a random number between 0 and the sum total

Use the random number + sum total to pick a random component

# Average Scores

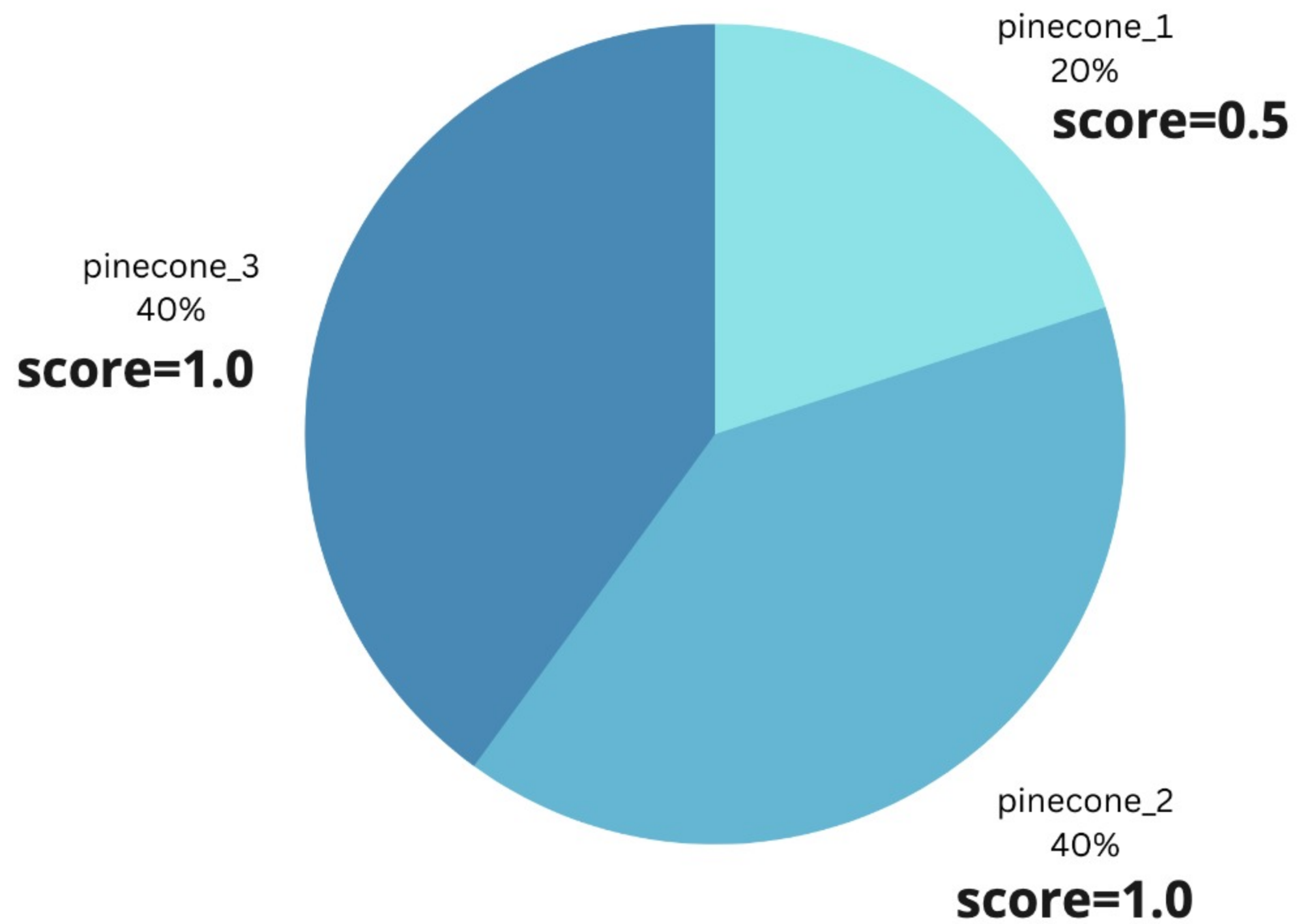| | |
|---|---|
| pinecone_1 | 0.5 |
| pinecone_2 | 1.0 |
| pinecone_3 | 1.0 |

Sum Total of Ratings → **2.5**

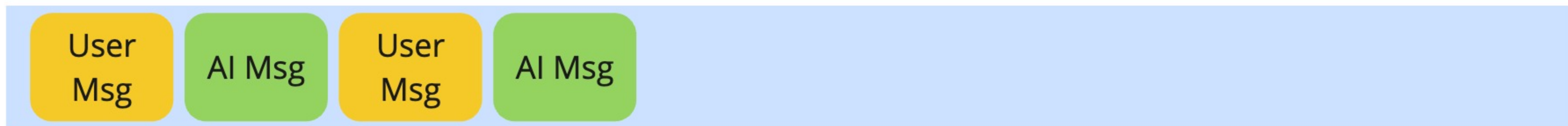Random number between 0 and that sum total → **1.3**

## ConversationBufferMemory

Stores all messages of
the conversation

## ConversationBufferWindowMemory

Stores only the last **K=2**
exchanges