

## برنام‌خدا

### مقدمه

سلام، اول از همه بخوام بگم به‌طور کلی پروژه چی بود: چت‌بات. کاربران با ثبت‌نام می‌تونن وارد سیستم شوند و با چت‌بات‌های موجود چت کنند. چت‌بات‌ها توسط یک سری کاربر خاص ساخته می‌شن که بهشون کاربر چت‌بات‌ساز می‌گیم. کاربر چت‌بات‌ساز هم در همین سامانه و در قسمت مخصوص خودش می‌تونه چت‌بات خودش رو بسازه. کاربر چت‌بات‌ساز توسط ادمین مشخص می‌شه. مهم‌ترین قسمت پروژه چت‌بات‌ها هستند. چت‌بات‌ها knowledge base هستند، به این صورت که کاربر چت‌بات‌ساز، چت‌بات خودش رو بر اساس دیتاهای مدنظرش می‌سازه. حالا کاربران می‌تونن با اون چت‌بات چت کنند. چت‌بات هم با استفاده از اون دیتاهایی که توسط کاربر چت‌بات‌ساز در اختیارش قرار گرفته شده، به سوالات کاربر پاسخ می‌ده.

این سامانه رو با استفاده از جنگو، [داکر](#)، [هم‌روش](#) (با استفاده از پایپ‌لاین [CI/CD](#) ساخته شده) و پست‌گرس پیاده‌سازی کردم. چت‌بات‌ها هم با استفاده از API مربوط به openAI طراحی شده‌اند.

### چت‌بات

همون‌طور که گفتم مهم‌ترین بخش کار، چت‌بات‌ها بودند. هر چت‌بات دیتاهای خاص خودش رو داره. که براساس اون دیتاها به سوالات کاربر پاسخ داده می‌شه. هر چت‌بات یک پرامپت سیستمی هم داره که توسط کاربر چت‌بات‌ساز مشخص می‌شه. به عنوان نمونه:

```
system_prompt="Based on the specific data: {data}, answer to the following question."
```

از اونجایی که هر چت‌بات می‌تونه دیتاهای خیلی زیادی داشته‌باشه، یک از چالش‌های مطرح این بود که کدوم بخش از این دیتاها به پرامپت کاربر مربوط هست. توجه داشته‌باشید که با توجه به محدودیت [توکن](#) موجود در پرامپت‌های openAI نمی‌تونستم کل دیتا رو برای openAI بفرستم. حتی اگر محدودیت توکن وجود نداشت، فرستادن کل دیتا هزینه‌ی الکی‌ای هست. حتی شاید فرستادن کل دیتا روی دقت پاسخ‌ها هم تأثیر بذاره.

**پیدا کردن دیتای مربوطه!**

برای پیدا کردن دیتای مربوط به پرامپت کاربر از [embedding](#) خود openAI استفاده کردم. اول از همه باید دیتای چت بات رو امبد می کردم. پرامپت کاربر رو هم امبد می کردم. شبیه ترین دیتای چت بات به پرامپت کاربر رو که پیدا کردم دیگه میتونم یک درخواست برای openAI بفرستم. اینجا هم یک سری مسائل به وجود میاد. اگر بخوام با هر پرامپت کاربر کل دیتاهام رو امبد کنم؛ هم خیلی طول می کشه، هم هزینه زیادی داره. و اینکه چجوری شبیه ترین دیتا به پرامپت کاربر رو پیدا کنم. در ادامه این ها رو بررسی می کنیم.

## ذخیره کردن و جست و جو در امبدینگ ها

زمانی که کاربر چت بات ساز یک دیتای جدید وارد می کنه یا یکی از دیتاهای قبلی رو ویرایش می کنه اون دیتا رو امبد می کنم. امبد دیتا رو داخل vector database ذخیره می کنم. حالا وقتی کاربر یک پرامپت می ده، اول پرامپت کاربر رو امبد می کنم تا با امبد های داخل وکتور دیتابیس مقایسه کنم و شبیه ترین دیتا رو پیدا کنم. برای جست و جو هم خود پستگرس قابلیتش رو داره.

## توهم ([hallucination](#))

یکی از مشکلات جدی مدل های زبانی hallucination هست یعنی وقتی یک چیزی رو بلد نیستن شروع میکنن به چرت و پرت جواب دادن. برای جلوگیری از این مشکل راه های زیادی هست. من به انتهای پیام کاربر (نه سیستم پرامپت!) متن زیر رو اضافه کردم.

```
usermessage += "\nSTRICTLY Do not give me any information about anything that is not mentioned in the PROVIDED CONTEXT."
```

## به خاطر آوردن پیام های قبلی

کاربران از یک چت بات انتظار این رو دارند که واقعا حالت چت داشته باشه و چت بات پیام های گذشته رو به خاطر داشته باشه. برای این کار من پرامپت کاربر رو به دو بخش تقسیم کردم. قسمت اول همون پیام حال حاضر کاربر که در قسمت های قبل بهش پرداختیم. قسمت دوم تاریخچه ی چت هست. برای ذخیره کردن تاریخچه ی چت از خود openAI استفاده کردم. بعد از اینکه چت بات به پرامپت کاربر پاسخ داد. پرامپت کاربر و پاسخ چت بات رو کنار هم قرار می دم و از مدل زبانی می خوام که یک خلاصه ی یک جمله ای از اون مطلب به من بده. خلاصه ای که از مدل زبانی گرفتم رو به تاریخچه ی جست و جو اضافه می کنم. در درخواست های بعدی تاریخچه ی چت رو در ابتدای پرامپت کاربر قرار می دم تا مدل زبانی پیام های قبلی رو به خاطر بیاره. سپس باز تاریخچه ی چت رو به روزرسانی می کنم. برای مطالعه ی بیشتر از این [لینک](#) استفاده کنید.

## دیس لایک

کاربر اگر پاسخ خوبی دریافت نکرد، بر روی علامت دیس لایک مقابل پاسخهای چت بات کلیک می کند. یک جواب دیگر مانند گذشته برای کاربر تهیه می کنم. تا به الان کار خاصی برای این قسمت انجام ندادم و فقط از دوباره به openAI ریکوئست می زنم به امید این که جواب بهتری تولید کنه. البته تغییراتی برای آینده در نظر دارم. یکی از راهها این است که به چت بات بگیم پاسخ به این سوال رو دوست نداشتم و یک جواب دیگه تولید کن. یا می توانیم از کاربر بپرسیم چرا از این پاسخ ناراضی بوده و این دیتا رو هم در اختیار چت بات قرار بدیم تا بتونه جواب بهتری تولید کنه.

## عنوان چت

برای اینکه چت ها عنوان داشته باشند؛ باز هم از خود openAI استفاده کردم. با دادن اولین پیام کاربر به چت بات ازش می خوام که برام یک عنوان تولید کنه.

## فول تکست سرچ

در قسمت چت کاربران می توانند در بین پیام های موجود جست و جو انجام بدهند. برای جست و جو از [full text search](#) پستگرس استفاده کردم. پستگرس یک سرچ وکتور فیلد داره که با استفاده از اون سرچ وکتور هر پیام رو داخلش ذخیره می کنم که لازم نباشه با هر بار جست و جو همش رو از اول بسازم. هنگامی که کاربر کوئری خود را وارد کرد برای آن یک کوئری سرچ می سازم و آن را با وکتورهای پیام ها مقایسه می کنم. و نتایج به دست آمده را بر اساس از رنک بالا به پایین نمایش می دم.