

Success Story: Liip | ZüriCityGPT

Mittels ChatGPT schnell an öffentliche Verwaltungsinformationen kommen

Die Suche nach Informationen auf Verwaltungswebseiten kann manchmal eine Herausforderung sein. Obwohl die Webseiten oftmals eine Fülle an Informationen bereithalten, sind sowohl die Struktur als auch die Suchfunktionen nicht immer optimal. Die Lösung? Das ChatGPT-Universum nutzen!

18.10.2023

OpenAI hat mit ChatGPT grosse Sprachmodelle (Large Language Model: LLM) populär und einer breiten Öffentlichkeit zugänglich gemacht. Dank benutzerfreundlichen Schnittstellen (APIs) und vernünftigen Preisen ist es nun einfach, eigene Experimente mit LLMs zu starten und individuelle Nutzungsszenarien zu realisieren. Und hier setzt unsere Vision an: Mit dieser Technologie einen einfacheren Zugang zu einer oft benötigten Informationsquelle schaffen. So entstand die Idee von ZüriCityGPT – ein freundlicher Bot, der nahezu alles weiss, was auf der Website der Stadt Zürich zu finden ist und (meistens) zuverlässige Antworten auf spezifische Fragen rund um die Verwaltung liefert.

Unser Ansatz

Um unser Ziel zu erreichen, haben wir uns für den technisch vernünftigsten Weg entschieden. Wir laden die gesamte Website mit einem Crawler herunter und zerlegen die Informationen in leicht verdauliche Einheiten für das LLM. Diese Segmente werden mittels sogenannter Embeddings vektorisiert und in einer Vektordatenbank gespeichert. Mit diesen Embeddings finden wir schnell die passenden Dokumente für eine gestellte Frage, um sie zur Beantwortung zusammen mit der ursprünglichen Frage an das LLM zu senden. Zum Schluss liefern wir die Antwort inklusive Quellenangaben direkt zurück. Diesen Ansatz nennt man Retrieval Augmented Generation (RAG).

Die Quellenangaben waren nicht nur für uns sehr wichtig, sondern erweisen sich auch als äusserst nützlich. Da ein LLM manchmal «halluziniert» und Dinge erfindet, helfen diese Angaben den Nutzenden, die Aussagen zu überprüfen. Wir arbeiten kontinuierlich daran, die Zuverlässigkeit zu verbessern, aber die perfekte Lösung haben wir (und andere) bisher noch nicht gefunden. Diese eigenwilligen Charakterzüge liegen in der Natur von LLMs.

Herausforderungen und Lösungen

Eine weitere Herausforderung ist es, die wirklich passenden Dokumente zu einer Frage zu finden. Die Prompt-Grösse (Anzahl der Zeichen, die dem LLM auf einmal geschickt werden können) ist beschränkt, und die Kosten steigen überproportional bei Modellen mit grösserem Kontext. Mit Experimentierfreude und aktivem Monitoring haben wir jedoch ein gutes Gleichgewicht gefunden. Dennoch werden weitere Analysen und Anpassungen folgen.

Auch ein Ansatz, der ausschliesslich auf semantischer Suche basiert, ist nicht immer ideal. Ein hybrides Modell, das auch eine klassische Suche mit einem invertierten Index und Stichwortsuche beinhaltet, könnte in bestimmten Fällen bessere Antworten liefern. Bisher haben wir diesen Ansatz aus Aufwandsgründen weitgehend vermieden, ausser bei Fragen zu Personennamen. Dank einer umfangreichen Liste von Vornamen können wir solche Anfragen erkennen und eine Volltextsuche in der Datenbank nach diesen Namen durchführen.

Eine beliebte Anfrage bei ZüriCityGPT ist die Suche nach Abfallabfuhr-Daten für bestimmte Stadtteile. Um diese besser beantworten zu können, benutzen wir nun die «Functions API» von OpenAI und leiten solche Fragen an die OpenERZ API weiter, was es uns ermöglicht, genaue und aktuelle Antworten zu geben, wie etwa wann die nächste Kartonsammlung an der Quellenstrasse sein wird. Eine weitere implementierte Anwendungsmöglichkeit war, Freizeit und Tourismus Fragen bei Zürich Tourismus abzufragen, da die Website der Stadt Zürich kaum solche Informationen hat.

Modell-Updates und Datenschutz

OpenAI bietet derzeit zwei Hauptmodelle an: GPT-3.5 und GPT-4. Zurzeit nutzen wir noch GPT-3.5 für ZüriCityGPT, da es viel kosteneffizienter und schneller ist, ohne viel an Zuverlässigkeit einzubüssen. GPT-4 zeigt jedoch vielversprechende Ergebnisse für gewisse Anwendungsfälle und halluziniert grundsätzlich einiges weniger. Wir sind jedenfalls gespannt auf die Entwicklungen in der Zukunft.

Beim Einsatz von gehosteten LLMs taucht schnell das Thema Datenschutz auf, insbesondere bei OpenAI, wo die Datenhandhabung nicht immer vollständig transparent ist. Die ohnehin öffentlich verfügbaren Daten werden nicht permanent bei OpenAI gespeichert, sondern nur auf unseren Servern. Wir nehmen den Datenschutz jedoch sehr ernst und prüfen auch alternative Hosting-Lösungen wie Azure OpenAI, die verbesserte Datenschutzrichtlinien und GDPR-Konformität bieten.

Für diejenigen, die eine Cloud-freie Lösung bevorzugen, könnte der Einsatz von Open-Source-LLMs eine Überlegung wert sein, obwohl dies aktuell noch mit erheblichen Initial- und Betriebskosten verbunden ist. Dies wird bei uns sicher immer mehr zum Thema; bis jetzt hatten wir aber die Ressourcen für eine vertiefte Auseinandersetzung damit noch nicht.

Ausblick

Die Reaktionen auf ZüriCityGPT waren überwältigend positiv, und wir freuen uns über das grosse Interesse und die Anfragen für weitere Versuche. In Kürze werden wir mehr über unsere Fortschritte und neue Projekte berichten können.

Als weiteres öffentliches Experiment haben wir LinkedDataGPT gestartet, mit den öffentlich verfügbaren Linked-Data-Daten der Stadt Zürich. Das Problem an Linked (und Open) Data ist normalerweise, dass diese schwer zugänglich sind, da es relativ viel Spezialwissen braucht, um die Daten abzufragen. Aber warum nicht die Möglichkeiten eines LLMs nutzen, um diese Abfragen zu schreiben? Gedacht, getan und mit einer einfachen Frage können nun die verfügbaren Daten abgefragt werden, wie etwa: «Welches sind die 20 häufigsten Vornamen?» oder: «Welche Religionen sind am beliebtesten?» Bei komplexeren Anfragen scheitert unsere Implementation zurzeit zwar noch, da besteht noch viel Raum für Optimierung. Das Interesse ist jedenfalls gross, diese Daten zugänglicher zu machen.

Fazit

ZüriCityGPT zeigt, wie ein individuell gestalteter Chatbot auf LLM-Basis mit relativ wenig Aufwand einen erheblichen Nutzen bringen kann, insbesondere bei der Verbesserung des Zugangs zu Informationen auf Behördenseiten. Die Möglichkeit, in mehreren Sprachen zu antworten, könnte die Barriere, für beispielsweise nicht deutschsprachige Personen, erheblich senken.

Darüber hinaus könnte eine LLM-basierte Suche auch intern bei Behörden einen echten Mehrwert bieten, indem sie das Finden von Informationen im Intranet erleichtert, obwohl dies natürlich mit eigenen Herausforderungen in Bezug auf Datenschutz und Zugriffsrechte einhergeht. Doch mit den technologischen Fortschritten, die wir heute sehen, scheint (fast) alles möglich.

Liip AG

Limmatstrasse 183 8005 Zürich contact@liip.ch 043 588 13 78 www.liip.ch

Liip ist eine Schweizer Digitalagentur mit starkem Mindset. Seit über zehn Jahren begleitet Liip Kundinnen und Kunden bei ihren digitalen Herausforderungen – von der Entwicklung preisgekrönter Webplattformen, Mobile Apps und Onlineshops bis zum Coaching für neue Arbeitsformen. Und das von A bis Z, mit Expertinnen und Experten für Strategie, Ideation, User Experience und Custom Development an Bord. Egal ob Start-up, Grossunternehmen oder Bundesbehörden, von Detailhandel bis Mobilität. Keine Lösungen von der Stange, sondern Digital Progress!