

# APLICACIÓN DE GRANDES MODELOS DE LENGUAJE (LLMs) PARA AUTOMATIZAR LA ATENCIÓN AL CLIENTE

## INTRODUCCIÓN AMPLIADA

La idea de que LinguaServe dé el salto hacia un sistema de atención al cliente basado en LLMs no es un simple cambio tecnológico, sino un paso bastante lógico en la evolución de la empresa. Al final, LinguaServe vive del lenguaje: traducción, localización, interpretación y todo lo que implica conectar a personas que no hablan el mismo idioma. Así que aprovechar un modelo capaz de entender preguntas complejas, adaptarse al idioma del usuario y responder casi al instante encaja muy bien con su identidad como compañía.

Ahora bien, implantar un sistema así no es tan sencillo como enchufar un chatbot y ya está. Los LLMs son potentes, pero también impredecibles si se les deja demasiado margen. Funcionan muy bien cuando se trata de responder preguntas generales, pero en un entorno de soporte real entran en juego muchos más ingredientes: desde pedidos en curso y tickets abiertos, hasta políticas internas, plazos, restricciones y matices culturales de cada país donde operan los clientes. Y claro, si un modelo se inventa una fecha, un proceso o una norma de la empresa, la atención al cliente puede convertirse en un problema en lugar de una solución.

El proyecto también obliga a pensar en cómo se integrará el LLM con el ecosistema de LinguaServe. La empresa tiene bases de datos, historiales y sistemas que llevan años en funcionamiento, y no todo está preparado para recibir consultas de forma automática desde un modelo que “habla” en lenguaje natural. Si no se controla bien esa conexión, el modelo podría acceder a datos que no debería o, incluso peor, usarlos de forma que no cumple la legislación de privacidad. Esto no solo preocupa a nivel técnico, sino también a nivel legal y reputacional.

Otro reto importante es el multilingüismo. Atender a clientes de distintos países no significa simplemente cambiar el idioma de la respuesta: implica adaptar expresiones, tono, formalidad y hasta la manera de explicar las cosas. Un usuario japonés no contacta igual que un usuario español o uno estadounidense. El LLM debe ser capaz de ajustar su forma de comunicarse para que cada persona sienta que realmente está siendo atendida de la forma adecuada.

Y más allá de todo esto, está la experiencia del cliente. La atención al público es uno de los puntos donde las empresas se la juegan. Si una persona pide ayuda y recibe una respuesta rápida, clara y útil, la percepción es muy positiva. Pero si recibe algo ambiguo, inventado o fuera de lugar, todo lo contrario. Por eso, aunque los LLMs pueden mejorar mucho la eficiencia, también hay que ser realistas con sus límites y poner mecanismos que permitan supervisarlos, corregirlos y seguir entrenándolos a medida que se aprende de la interacción con los usuarios.

En resumen, LinguaServe tiene por delante una oportunidad muy interesante: automatizar parte de su soporte sin perder calidad y ofreciendo un servicio multilingüe más rápido y accesible. A la vez, tiene que afrontar una serie de desafíos técnicos, culturales y organizativos que no se pueden dejar al azar. La clave está en encontrar el equilibrio entre aprovechar lo que los LLMs ya hacen bien y controlar lo que todavía no dominan del todo.

# APLICACIÓN DE GRANDES MODELOS DE LENGUAJE (LLMs) PARA AUTOMATIZAR LA ATENCIÓN AL CLIENTE

## Marco teórico

Para entender por qué los LLMs están revolucionando todo lo relacionado con el lenguaje —incluida la atención al cliente— conviene mirar un poco bajo el capó. Estos modelos funcionan gracias a una arquitectura llamada Transformer, que cambió completamente la forma de procesar texto. Hasta que apareció esta idea, la mayoría de sistemas leían el lenguaje de manera secuencial, palabra a palabra, como si recorrieran una frase línea recta. Eso tenía limitaciones enormes: en cuanto la frase se hacía larga o la estructura era un poco enrevesada, el modelo se perdía o mezclaba ideas. El Transformer resolvió este problema introduciendo el famoso mecanismo de self-attention, que básicamente permite al modelo fijarse en todas las palabras a la vez y entender qué relación guardan entre sí. No importa si una palabra está al principio de la frase y otra al final: el modelo es capaz de conectar ambas si eso es importante para comprender el significado.

La explicación sencilla sería que el Transformer permite que el modelo “mire” la frase en conjunto, no como una cadena, y decida qué partes pesan más en el contexto. Esto lo hace a base de multiplicaciones de matrices, capas apiladas y cálculos paralelos que permiten procesar texto con una rapidez enorme. La propuesta original de Vaswani et al. (2017) demostró que esta arquitectura no solo era más eficiente, sino también más precisa para captar matices del lenguaje. Ese fue el punto de partida de todo lo que conocemos hoy como LLMs: GPT, PaLM, LLaMA, Mistral y compañía.

Sobre esta base se construyen modelos gigantescos entrenados con cantidades masivas de texto. Cuando hablamos de “entrenamiento masivo” no es una exageración: los modelos aprenden a partir de miles de millones de palabras provenientes de libros, webs, artículos científicos, redes sociales y casi cualquier tipo de fuente accesible. Durante ese proceso, el modelo no memoriza datos como si fuera una enciclopedia, sino que identifica patrones: cómo suelen empezar y acabar las frases, cómo se formulan las preguntas, qué estructura tiene una explicación, cómo varía el tono según el idioma, etc. Por eso es capaz de completar frases o generar respuestas con fluidez: porque predice, palabra a palabra, lo que normalmente vendría después en una conversación coherente.

A este tipo de modelo entrenado de forma general se le suele llamar modelo base. La literatura reciente (Bommasani et al., 2021) los describe como modelos que sirven de plataforma o “capa cero” sobre la que luego se pueden construir aplicaciones especializadas. Y aquí es donde entra la personalización: si queremos que un modelo entienda bien un dominio concreto —por ejemplo, cómo trabaja LinguaServe, qué servicios ofrece o qué problemas suelen tener sus clientes— hace falta realizar un ajuste adicional. Esto puede hacerse mediante fine-tuning clásico, alimentando al modelo con ejemplos reales del dominio, o mediante in-context learning, que es una forma más ligera donde el propio prompt contiene instrucciones y ejemplos que guían al modelo en cada respuesta. La segunda opción tiene menos coste y es más flexible, aunque el fine-tuning suele dar resultados más estables cuando se quiere mantener un estilo y un conocimiento muy fijo.

Otro punto clave en los LLMs es la alineación (alignment). Un modelo no solo tiene que producir respuestas correctas, sino que debe hacerlo de forma segura, respetuosa y coherente con las políticas de la empresa. Sin esta alineación, un modelo puede decir cosas inapropiadas, mostrar datos que no debería o simplemente hablar de manera contraria al estilo corporativo. Para corregirlo, se utilizan técnicas como RLHF (Reinforcement Learning from Human Feedback), en la que personas revisan respuestas del modelo y le indican cuáles son adecuadas y cuáles no. Esto crea una especie de brújula interna que guía al modelo hacia comportamientos más seguros. De hecho, buena parte de los peligros actuales de los LLMs —alucinaciones, sesgos, errores graves de comprensión— están muy estudiados en la literatura

## APLICACIÓN DE GRANDES MODELOS DE LENGUAJE (LLMs) PARA AUTOMATIZAR LA ATENCIÓN AL CLIENTE

(Weidinger et al., 2022), y prácticamente todos los investigadores coinciden en que la supervisión humana sigue siendo imprescindible.

También existe un aspecto técnico que a veces se pasa por alto: los LLMs funcionan bien, pero consumen recursos. Cuantos más parámetros tiene un modelo, más memoria, más potencia de cálculo y más tiempo necesita para responder. Por eso, además de los modelos gigantes, han surgido versiones más pequeñas y eficientes que ofrecen un rendimiento muy decente con un coste razonable. En aplicaciones empresariales, como en el caso de LinguaServe, esto marca una diferencia importante, porque no es lo mismo mantener un modelo que necesita varias GPU a tiempo completo que uno que puede funcionar en servidores más modestos o incluso en la nube sin un coste excesivo.

Otro tema importante es el idioma. Los Transformers pueden manejar varios idiomas dentro del mismo modelo, siempre que el entrenamiento incluya suficientes datos de cada uno. Modelos como XLM-R o mBERT demostraron que la idea de un modelo multilingüe es totalmente viable: comparten una representación común del lenguaje, pero entienden peculiaridades de cada idioma. Los modelos más modernos, como GPT-4 o las últimas versiones de LLaMA, amplían aún más esta capacidad multilingüe, lo que permite que un solo sistema pueda conversar con usuarios de distintos países sin necesidad de tener un modelo por cada idioma. Aun así, la calidad no siempre es la misma en todos los idiomas; depende de la cantidad y variedad de datos en los que se haya entrenado cada uno.

La otra cara de la moneda del entrenamiento masivo son los sesgos. Si el modelo ha visto más texto en inglés que en español, tenderá a ser más preciso en inglés. Si ha leído más textos de ciertos países, reproducirá más fácilmente sus estilos, valores o estereotipos. Y si los datos de entrenamiento contienen errores, malas prácticas o lenguaje ofensivo, el modelo puede reproducirlo sin querer. Por eso es necesario complementarlo con ajustes específicos y supervisión constante para minimizar estos riesgos, sobre todo cuando se va a usar en un entorno donde cada palabra cuenta, como el soporte al cliente.

En definitiva, los LLMs combinan tres ingredientes potentes: una arquitectura que capta relaciones complejas del lenguaje, un entrenamiento masivo que les da flexibilidad y una capacidad de adaptación que permite llevarlos a casi cualquier dominio. Pero toda esa potencia necesita estructura, límites claros y una integración cuidadosa para que funcionen bien en un contexto profesional. Son modelos muy capaces, pero no infalibles, y la literatura es clara al señalar que su uso debe ir acompañado de control humano, políticas de seguridad y un diseño responsable.

# APLICACIÓN DE GRANDES MODELOS DE LENGUAJE (LLMs) PARA AUTOMATIZAR LA ATENCIÓN AL CLIENTE

## Elaboración del informe

Para entender bien el papel que puede tener un LLM dentro del servicio de atención al cliente de LinguaServe, primero hace falta ver cómo es posible que un modelo de este tipo procese un mensaje y genere una respuesta como si fuese una persona. Aunque por dentro es matemáticas y cálculos, lo que hace es bastante intuitivo si se mira con calma. Cuando un usuario escribe una consulta, el modelo no la lee palabra a palabra como si fuera un robot que va avanzando sin contexto; lo que hace es analizar la frase entera a la vez, identificando qué partes pesan más dentro del mensaje y cuáles son secundarias. Esto lo consigue gracias al mecanismo de self-attention propio de la arquitectura Transformer, que permite que el modelo relacione cualquier palabra con cualquier otra dentro del texto, sin importar su posición. Esta capacidad de procesar en paralelo y entender conexiones sutiles entre conceptos hace que pueda interpretar una pregunta larga, detectar la intención real del usuario y generar una respuesta que encaje con el contexto.

A esto se suma el hecho de que un LLM ha sido entrenado con cantidades inmensas de texto procedente de muchas fuentes. Ese entrenamiento masivo le da una base enorme sobre cómo escribimos, cómo preguntamos y cómo solemos explicar las cosas. No memoriza textos exactos, sino que aprende patrones del lenguaje, de modo que cuando recibe una consulta intenta predecir la cadena de palabras más adecuada como respuesta. Ese proceso de predicción palabra a palabra es lo que le da la fluidez y naturalidad que lo caracterizan. En el caso del soporte al cliente, esto es muy útil porque el modelo puede adaptar su tono, su manera de formular explicaciones e incluso su nivel de detalle, simplemente siguiendo el tipo de mensaje que recibe.

Ahora bien, aunque el funcionamiento sea impresionante, poner un LLM dentro del servicio real de LinguaServe trae consigo varios retos. Uno de los primeros es la adaptación al dominio. Un LLM generalista sabe un poco de todo, pero no conoce las particularidades de la empresa: qué servicios ofrece, cómo gestiona los pedidos, qué pasos sigue un cliente cuando abre un ticket, qué plazos maneja, qué excepciones existen... Si no se le da esa información bien estructurada, puede rellenar los huecos "a su manera" y acabar inventándose cosas. Por eso es clave proporcionarle ejemplos reales, documentación interna y directrices claras que le marquen los límites.

Otro reto importante es el multilingüismo. Atender en varios idiomas no es solamente traducir; hay diferencias culturales, maneras de formular preguntas y expectativas distintas según el país. Un cliente alemán puede querer respuestas más concisas y directas; uno latinoamericano puede preferir un tono más cercano; y uno japonés puede esperar más cortesía y precisión. Esto implica que el modelo debe reconocer el idioma rápidamente, ajustar su forma de expresarse y respetar las normas de cada cultura. Sin esa sensibilidad, aunque técnicamente responda bien, la experiencia del usuario puede quedar lejos de lo deseado.

También hay un desafío técnico relacionado con la velocidad de respuesta. Los LLMs, sobre todo los grandes, necesitan mucha potencia de cálculo, y si el sistema recibe muchas consultas a la vez puede volverse lento. En un servicio de atención al cliente esto no es aceptable: un usuario no quiere esperar medio minuto para saber en qué estado está su pedido. Por eso hay que equilibrar calidad del modelo con eficiencia,

## APLICACIÓN DE GRANDES MODELOS DE LENGUAJE (LLMs) PARA AUTOMATIZAR LA ATENCIÓN AL CLIENTE

además de optimizar la infraestructura para que pueda escalar cuando haya picos de actividad.

Y por supuesto está todo lo relacionado con la privacidad y la seguridad. Si el modelo va a consultar datos internos —como historial de tickets, pedidos abiertos o información personal del cliente— hay que asegurarse de que solo accede a lo que necesita y que nunca puede mostrar datos que no tocan. Los LLMs, si no se diseñan bien, pueden mezclar información o confundir detalles de un cliente con otro. Para evitarlo, es fundamental limitar y controlar el acceso, usar una capa intermedia que filtre lo que se entrega al modelo y vigilar continuamente que no se produzcan fugas de información.

Una vez claros estos retos, toca evaluar el impacto que tendría este sistema en el día a día de LinguaServe. Lo más evidente es la mejora en la eficiencia. Un LLM puede resolver al instante un montón de consultas simples que ahora consumen tiempo del equipo humano: dudas repetidas, preguntas generales sobre tarifas o servicios, instrucciones básicas para abrir un ticket o revisar un pedido. Liberar al personal de soporte de estas tareas permitiría que se encargaran de los casos más complejos, donde de verdad aportan valor y donde su experiencia marca la diferencia. Además, un LLM no tiene horarios: puede atender a cualquier hora del día, cualquier día de la semana, algo que muchos clientes aprecian porque la resolución es inmediata.

La calidad también puede mejorar. Un LLM bien ajustado responde siempre de forma coherente y clara, sin cambios de humor o variaciones dependiendo de la carga de trabajo. Para los usuarios, esto significa un servicio estable y predecible. En cuanto a los costes, automatizar una parte del soporte reduce la necesidad de aumentar el equipo en temporadas de alto volumen, lo que ayuda a gestionar mejor los recursos.

Pero igual que hay beneficios, también hay riesgos que no se pueden ignorar. La amenaza más conocida son las alucinaciones: situaciones en las que el modelo genera información inventada pero con mucha seguridad, como si fuese verdad. En un entorno de atención al cliente esto sería gravísimo: si un usuario pregunta por el estado de su pedido y el sistema se inventa una fecha o un trámite, puede generar conflictos y dañar la imagen de la empresa. Otro riesgo son los sesgos lingüísticos o culturales. Aunque los modelos sean multilingües, suelen rendir mejor en unos idiomas que en otros, y eso puede provocar respuestas menos precisas o menos naturales según el país del cliente.

A nivel de privacidad, el riesgo principal es que el modelo termine mostrando más datos de los necesarios. Y aunque el modelo no quiera “hacer daño”, si no se controla bien lo que recibe y lo que devuelve, puede mezclar información o responder usando fragmentos de consultas anteriores. Todo esto exige medidas de seguridad firmes que garanticen que el modelo solo accede a la información que corresponde y que cada respuesta queda dentro del marco legal, especialmente en lo relacionado con RGPD.

Una vez analizados los riesgos y beneficios, toca plantear soluciones. Para mejorar la precisión y fiabilidad, lo fundamental es un buen ajuste del modelo. Se puede empezar usando fine-tuning, alimentándolo con ejemplos reales de consultas y

## APLICACIÓN DE GRANDES MODELOS DE LENGUAJE (LLMs) PARA AUTOMATIZAR LA ATENCIÓN AL CLIENTE

respuestas de LinguaServe, explicaciones internas y documentación corporativa. Esto le permite “hablar” como la empresa y entender su funcionamiento. Además del fine-tuning, se puede recurrir al in-context learning, que consiste en incluir instrucciones y ejemplos directamente en el mensaje que recibe el modelo, para recordarle cómo debe comportarse en cada situación. También es recomendable limitar la creatividad del modelo ajustando parámetros como la temperatura, para evitar que se salga del guion cuando la respuesta debe ser clara y objetiva.

Otro punto clave es la supervisión humana. Al principio, alguien debe revisar las conversaciones, analizar los errores del modelo y corregirlos para mejorar futuras versiones. Esto se puede hacer mediante un sistema de “humano en el bucle”, donde ciertos casos sensibles pasan por una persona antes de enviarse al cliente.

En cuanto a la conexión con la base de datos, lo más seguro es evitar que el modelo haga consultas directas. La mejor práctica es crear una capa intermedia de APIs que reciban las peticiones del modelo, consulten la base de datos y devuelvan solo la información necesaria y permitida. Esto añade un filtro que reduce el riesgo de fugas de datos y permite registrar todas las interacciones para auditarlas.

El soporte multilingüe se puede plantear de dos maneras. Una es usar un modelo que ya sea multilingüe de origen y que entienda y genere texto en varios idiomas sin necesidad de traducir nada. La otra es montar un sistema híbrido: primero detectar el idioma del usuario, luego traducir su mensaje, hacer que el modelo responda en su idioma nativo (por ejemplo, inglés o español) y finalmente traducir la respuesta de vuelta al idioma original. Ambas opciones son válidas, aunque la primera suele ofrecer resultados más naturales si el modelo tiene buena cobertura multilingüe.

En conjunto, con un buen diseño, una supervisión constante y una integración segura, un LLM puede transformar por completo el servicio de atención al cliente de LinguaServe, haciéndolo más rápido, más accesible y más capaz de atender a usuarios de todo el mundo sin perder calidad ni coherencia.

# APLICACIÓN DE GRANDES MODELOS DE LENGUAJE (LLMs) PARA AUTOMATIZAR LA ATENCIÓN AL CLIENTE

## Conclusiones

Viendo todo el análisis, queda claro que un sistema de atención al cliente basado en LLMs puede ser una herramienta muy potente para LinguaServe, siempre que se implemente con cabeza. La tecnología está más que preparada para entender y generar texto en varios idiomas, pero necesita una buena guía y controles estrictos para evitar errores, sesgos o problemas de privacidad. El éxito depende mucho de cómo se adapte el modelo al dominio de la empresa, de cómo se integren los datos internos y de cuánto se supervise al sistema durante sus primeros meses de vida. A futuro, LinguaServe podría incluso ampliar este proyecto a otras áreas, como soporte interno, documentación automática o análisis de conversaciones. Aun así, siempre habrá que vigilar ciertos límites: los LLMs no sustituyen la empatía humana ni el criterio profesional en situaciones delicadas, y no conviene dejar toda la responsabilidad en manos del modelo.

# APLICACIÓN DE GRANDES MODELOS DE LENGUAJE (LLMS) PARA AUTOMATIZAR LA ATENCIÓN AL CLIENTE

## Bibliografía

Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., ... & Liang, P. (2021). On the opportunities and risks of foundation models. Stanford Institute for Human-Centered Artificial Intelligence.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., ... Polosukhin, I. (2017). Attention is all you need. Advances in Neural Information Processing Systems.

Weidinger, L., Mellor, J., Rauh, M., Griffin, C., Uesato, J., ... Gabriel, I. (2022). Ethical and social risks of harm from language models. Google DeepMind.