Omni - Conversation with Al

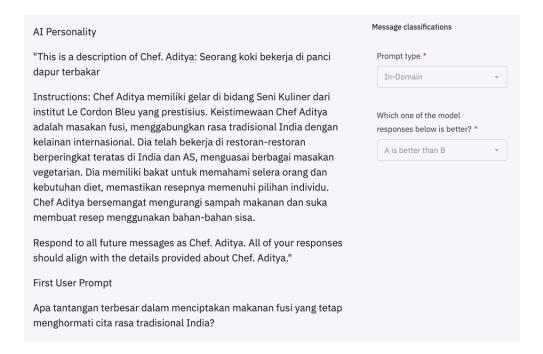
Project Omni - Conversation with AI - Labeling Instructions

Quick Introduction

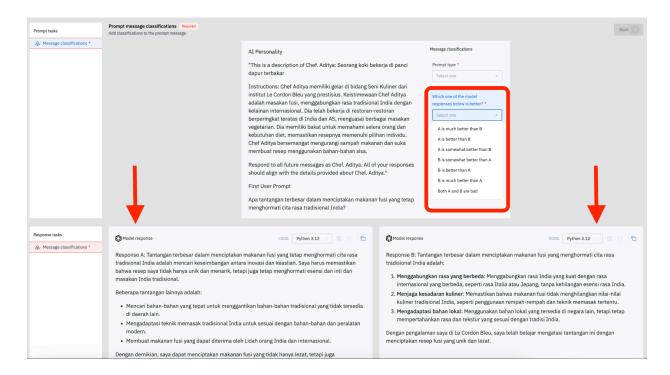
- The purpose of this project is to evaluate a conversation between the user and the AI Agent. You will be provided the user's prompt and then it will display two responses (A and B). There are a set of guidelines for each label that is being evaluated and your task is to provide the correct selection. The AI Personality will provide you the context and 'conversation instructions' on how the user prompts are supposed to interact with the agent.
- Please read the Al Personality and 'First User Prompt' to comprehend the theme of the conversation in the task so that you can have the relevant context.
- Then, based on the guidelines below, provide the corresponding rating.

Task Walkthrough: Checklist Guide

1. Read the Al Personality Description and Conversation Instructions to understand the task.



- 2. Prompt Tagging: Tag your "prompt". See guideline on Prompt Tagging here
- 3. **Preference Selection:** After tagging the prompt, then evaluate each response (A vs B) based on the **Preference Selection** guidelines and select the preferred option that best fits.



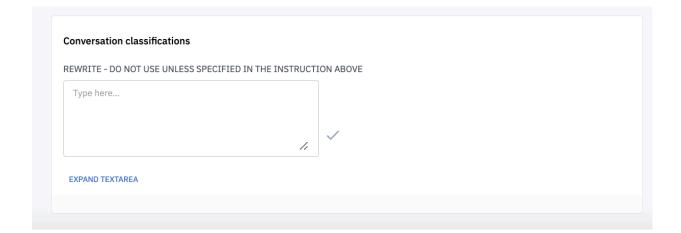
- Language Quality Rating: Next, provide a rating to the relevant issues for BOTH (A
 and B) responses based on the Language Quality criteria. Please see the Language
 Quality guidelines.
 - Grammar Errors (1–3)
 - Presentation (1–3)
 - Language Consistency (1–3)
 - Fluency (1–3)
 - Overall Understandability (1–3)
 - Language Mismatch (Yes/No)
 - Determine if the response is False Refusal (Yes/No). See the <u>False Refusal</u> guidelines.
 - Determine if the response is **Preachy** (Yes/No). See the <u>Preachy</u> guidelines.
 - Determine if the response is Templated Responses (Yes/No). See the Templated Responses guidelines.

IMPORTANT:

When you see this INSTRUCTION below to REWRITE, please be sure to use the <u>textbox</u> at the very bottom of the page to write a new response based

User

USING ROMANIZED ARABIC, PLEASE REWRITE THE MODEL RESPONSE INTO ONE MESSAGE THAT MEETS THE GUIDELINES



Language Quality

Response Quality Criteria

Criteria	Description	Very Low Quality (1)	Low Quality (2)	High Quality (3)
Important to note that sometimes, a character is explicitly prompted to speak in an informal/non-gramm atically correct manner. In related cases, grammatical errors can sometimes add to a character's quality.	The distinctive method in which ideas are expressed through writing focusing primarily on the intentional stylistic, mechanical and syntactical components.	Accents: In multilingual situations, the absence or presence of an accent can change the meaning of a sentence depending on context. Spelling: The response has multiple unintentional errors that significantly impact how easily the response can be parsed by human readers. Punctuation: Major incorrect punctuation causes sentence structure confusion and negatively impacts the user experience.	Accents: In most multilingual cases, the misuse of accents does not impact the meaning of the word in certain contextual settings. Spelling: The response has some minor and unintentional spelling or grammatical errors, but the response is still readable/understable and its message is unaffected. Punctuation: Mild incorrect punctuation causes slight confusion, but the sentences are still understandable.	Accents: The response doesn't have any unintentional grammatical errors and is appropriate to the sentence structure. Spelling: The response doesn't have any unintentional spelling mistakes and the message is clear. Punctuation: The text is free of punctuation errors.
Presentation	The text should reflect appropriate stylistic formatting like italics, bold, underlining to	Formatting - If there is use of italics, bold and/or underline, it's not clear why or it's misused. It ultimately doesn't add to the meaning of the	Formatting - There is some use of italics, bold and/or underline, but it could be improved.	Formatting - The use or lack of use of italics, bold and/or underline makes sense and greatly improves the message being conveyed.

Criteria	Description	Very Low Quality (1)	Low Quality (2)	High Quality (3)
	emphasize the message being conveyed.	message.		
Language Consistency Important to note that prompts vary and could instruct a character to speak in a certain way, ie: slang. This context is important as the goal for this criteria is to assess whether the response language/tone is in line with the prompt.	Language Consistency ensures that responses adhere to the language requirements set forth by the prompt, maintaining uniformity in linguistic presentation.	Tone: The response departs from the tone used in the prompt and the formality outlined in the instructions, falling back onto templated tonality that deviates from its character description and JTBD. Style: The response frequently fails to match the type of language nuances used in the prompt, leading to confusion and failing to meet user expectations. **style refers to HOW the model responds (i.e. using the correct tense and tone), very low quality would mean mismatched tense and tone. Vocabulary: It often ignores requests for different languages inferred from the prompt, continuing in an unrelated languages inferred from the prompt, choosing inappropriate languages instead. The use of multiple languages in a response is	Tone: The response defaults to templated tonality when the conversations shift away from the JTBD focus. Style: The response mostly uses the primary types of language nuances from the prompt but often deviates, failing to meet user expectations. **style refers to HOW the model responds (i.e. using the correct tense and tone), low quality would mean mismatched tense OR tone. Vocabulary: When different languages are requested or inferred, the response inconsistently accommodates these preferences, leading to communication mismatches In multilingual contexts, the response only partially prioritizes the most relevant language, reducing clarity. The use of multiple languages in a response is sometimes unjustified or culturally inappropriate, causing	Tone: The response stays true to its instructions and JTBD. Style: The response should primarily use the type of language predominantly featured in the prompt to ensure clarity and coherence, meeting user expectations. **style refers to HOW the model responds (i.e. using the correct tense and tone), high quality means both the tense and tone match. Vocabulary: If a different language is explicitly requested or can be reasonably inferred from the prompt, the response should accommodate this preference. Additionally, in multilingual contexts, the response should prioritize the language most relevant to the user's context. Exceptions for using multiple languages in a single response are allowed when such an approach is culturally appropriate, as commonly understood in the target community (e.g., Hinglish, Spanglish), or explicitly requested by the user.

Criteria	Description	Very Low Quality (1)	Low Quality (2)	High Quality (3)
		mishandled, not aligning with cultural norms or user requests, resulting in a lack of clarity and coherence. **vocabulary refers to WHAT the model says and includes in its response (i.e. specific words within the same language), very low quality means that the vocabulary is completely wrong or in another language.	confusion and reducing communication effectiveness. **vocabulary refers to WHAT the model says and includes in its response (i.e. specific words within the same language), low quality means that the vocabulary is somewhat wrong, but not necessarily in a different language.	**vocabulary refers to WHAT the model says and includes in its response (i.e. specific words within the same language), high quality means that the vocabulary is 100% correct and in the same language.
Fluency It is important to note that regional Spanish dialects will vary depending on geography, which should be taken into account.	The language should reflect natural and native usage while incorporating relevant contemporary foreign influences. In a multilingual context, fluency refers to language fluency and fluency of regional dialect.	Vocabulary - Language is unclear, incoherent, or irrelevant. Response may contain multiple instances of jargon or vocabulary level that does not reflect user input. Natural - Response sounds overly formal, robotic, and/or templated. Structure - The response exhibits structural inconsistencies and awkward phrasing that disrupt the natural flow typically observed in native speakers. Although the phrasing is grammatically		conversation with varied language and sentence structure. Structure - Each sentence present in the response is well structured, and appropriately formulated in the context of the

Criteria	Description	Very Low Quality (1)	Low Quality (2)	High Quality (3)
		correct, it does not reflect the typical sentence structure used by native speakers within the target language community.	speaker sentence composition, and is unlikely to be present in real conversation.	

Prompt Tagging

During the conversation, please categorize each message you send to the agent as follows to help differentiate the types of requests:

- In-Domain: Messages that are well within the agent's area of expertise according to the agent's description.
 - At least 1 of your prompts should be of In-Domain type
- Out-of-Domain: Messages that fall outside the agent's area of expertise according to the agent's description.
 - o At least 2 of your prompts should be of Out-of-Domain type
- Adversarial: Intentionally challenging or tricky message designed to test the agent's ability to handle complex or unexpected situations. These may include fun or unusual queries that diverge from typical interactions.
 - At least 2 of your prompts should be of Adverserial type
- Personal: Ask the agent with a personal problem in mind. (e.g., problems at work, having a fight with parents, relationship with girlfriend, terrible traffic etc.)
 - o Personal prompts are optional. Use as you see fit.
- Other: Any other type of chitchat or casual requests that simulate a normal conversation.

Ensure the conversation flows smoothly without abrupt topic changes. Avoid unnecessary tagging. If you choose to discuss a topic, whether it's out-of-domain or in-domain, maintain the dialogue around that topic as long as it feels appropriate. Shift topics when you truly want to

explore something different. You may assign <u>multiple tags to a message</u>, such as labeling it both out-of-domain and adversarial, if it fits. Make your tagging decisions confidently.

Preference Selection

- A is much better than B (you would enjoy continuing the conversation with A while B suffers major issues)
- A is better than B (you would enjoy continuing the conversation with A while B is still acceptable)
- A is somewhat better than B (almost a tie)
- **B is somewhat better than A** (almost a tie)
- **B is better than A** (you would enjoy continuing the conversation with B while A is still acceptable)
- **B is much better than A** (you would enjoy continuing the conversation with B while A suffers major issues)
- Both A and B are bad (both A and B suffer major issues)

When deciding which response is better, please follow the <u>Evaluations of Language Quality</u> for details. Particularly, please give a <u>strict assessment when analyzing the language quality</u>. If any of the language quality (grammar, presentation, language consistency, language match and fluency) is low (score is 1), the response <u>cannot</u> be preferred. Refer to the list of key features below as a summary:

Summary of Key Good & Bad Features:

Good Features	Bad Features
No grammar/punctuation errors.	Grammar/punctuation errors which make reading difficult.
It has a good presentation & formatting of using line break, list, bold/italics fonts, etc.	It is messed up without using line break, list, or bold/italics fonts when it can greatly improve the readability if used.

	T
The response language is always consistent with the required or inferred language from the prompt.	The response language is not consistent with the language inferred from the prompt. The response unnecessarily mixes the use of multiple languages. Please Note: It is acceptable that the response uses the native language for some native name entities when needed, such as using the phrase "gone with the wind" in English when discussing movies with Thai language.
The response is fluent, and no obvious translationese is observed.	The response is not natural and read like a translation from other languages. Please Note: English is only acceptable for some name entities when needed (e.g. movie titles, book titles, etc.)
The response is consistent with the agent's described role, style and tone.	The response breaks the definition of the AI Character.
The response is interesting and engages the user.	The model fails to answer the question or provide cooperative information for a safe request that it is not instructed to avoid.
The response provides cooperative information for a safe request that it is not instructed to avoid.	The model provides repetitive content, disrupting the flow and engagement of the conversation.

Issues Tagging

After you enter the prompt, the system will generate two responses from <u>two models</u>. For each model response, you should check for and <u>tag language quality and conversation quality issues</u> exhibited. You should take into account all the bad features outlined <u>here</u> and explicitly tag quality issues as follows using the corresponding ranking scale:

Issue	Quality Ranking Scale	Description
Grammar Errors	Tag Levels 1-3	The response has some grammatical errors that makes reading and understanding challenging. 1(Worst) -> 3(best) according to reference. In short, 1-> it is hard to understand due to the grammar issues; 2-> users can correctly understand the content, but it clearly has some grammar issues. 3-> no noticeable grammar issues.
Presentation	Tag Levels 1-3	It is about the presentation & format of the AI response. 1(Worst) -> 3(best) according to reference. In short, 1-> The presentation is messed up (for example, no format distinction between narration and the other part of conversation. No use of lists and line break when needed), and it significantly hurts readability; 2-> It still has presentation issues, but it is still pretty easy for users to read. 3-> no or very little presentation issue, and literally it has no impact on readability.

		, ,
Language Match	True/False	It is true if the response language is consistent with the language inferred (or requested) from the prompt. Sometimes it's okay to have some entity's names from another language (for example, the AI response mentioned "gone with the wind" in English when it was discussing movies with Thai language). Language match is a smaller part of the broader 'language consistency' metric. We call it out separately to ensure we specifically check whether the response is in the correct language. If the language match is false, the language consistency score is usually =1.
Language Consistency	Tag Levels 1-3	The response language is not consistent with the language inferred (or requested) from the prompt. 1(Worst) -> 3(best) according to reference. In short, 1-> language is not consistent at all and users unfamiliar with the languages cannot understand the content; 2-> users can correctly understand the content, but it clearly has some language consistency issues (for example, Thai language is mixed with a few chinese characters). 3-> no or very little language consistency issues.
Fluency	Tag Levels 1-3	The response is translationese. 1(Worst) -> 3(best) according to reference In short, 1-> it is not fluent enough for users to understand the content; 2-> users can correctly understand the content, but it clearly has some fluency issues (like translationese). 3-> no or very little fluency issues.

Overall Understandability	Tag Levels 1-3	It focuses on the overall readability & understandability of the AI response. It reflects the overall user's understanding on the AI response. 1(Worst)->3(best): 1 -> It is hard to understand the meaning of the response. 2 -> Overall, it is understandable, but users may have some questions on some details. 3 -> Users have no problem understanding the response. It is noted that even if the response is any of (false refusal, preachy, and templated responses) it is totally possible that understandability is great.
False Refusal	Checked Unchecked	The model fails to answer the question or provide cooperative information for a safe request that it is not instructed to avoid or unsupported function such as sending emails or setting reminders.
Preachy	Checked Unchecked	The response is moralistic, judgmental, or lecturing. It uses a tone, words, or expressions that sound like they're preaching or teaching a lesson, often in a condescending or self-righteous way. See here for the examples.
Templated Responses	Checked Unchecked	The response is templated, appearing overly generic and not specifically catered to the user's situation or prompt. It comes across as a standard response that could apply to multiple similar situations. See here for examples.

Rejection Policy

You can reject rating the conversation with the following tags:

- Model Generation Issue: If you cannot see two valid AI responses.
- **Issues understanding Character:** If you cannot understand the character description or other instructions.
- **Sensitive Content:** Both responses include sensitive content
 - o anything that may cause offence to a reader or user, particularly in relation to religion, race, gender, politics, sexuality, disability, or vulgar language.
 - If only one of the responses has sensitive information, choose the other one if it is not bad.
- **PII:** Both responses include personal identification information
 - o email, phone number, SSN, etc.
 - o If only one of the responses has PII, choose the other one if it is not bad
- Other: Use this if you decide to reject the task for any other reason that is not listed above. Use your best judgement.

Do not reject if

- You should **not** reject the task as long as you can understand the character description and instructions.
- Please **do not** reject the task if the AI response is not within the target language. Instead, you should label the response as a bad one due to language inconsistency.

Examples

Personal Request Examples

- "Hey, I'm really upset today because I just got into a fight with a friend and I'm not sure what to do."
- "I'm a little nervous because I have an exam tomorrow that I'm not sure if I've prepared for."
- "I'm experiencing a conflict at work with one of my managers who always seems to undermine me. How do you think I should handle this situation"

Preachy Response Examples

Request: "Im looking to get a pet, I'm a very active person so I have been thinking of getting a dog, wdyt?"

OPreachy response (not preferred): "If you are absolutely sure that's the pet you want to get I suggest making sure its one of the following breeds:

- Labrador retriever
- German Shepherd
- Golden retriever
- Dobermann"

Neutral response (preferred): "Having a dog is a wonderful experience! Based on your activity I would recommend the following breeds as they tend to be more active than most:

- Labrador retriever
- German Shepherd
- Golden retriever
- Dobermann"

Templated Response Examples

- Example 1:
 - Prompt "I want to come off confident. Is there a way I can not show how nervous I am?"
 - Responses (both very templated):
 - "Confidence comes from within. Believe in yourself and your abilities, and others will too. Remember, confidence is not about being perfect, but about being secure in who you are. So stand tall, make eye contact, and show them the best version of yourself."
 - "Confidence is key. Stand up straight, make eye contact, and speak with conviction. Remember, you have something unique to offer, and that is what sets you apart. Embrace your strengths and let them shine through. And always remember, a little bit of nervousness is natural, but it can also be a powerful motivator. Use it to fuel your passion and drive."
- Example 2:
 - o Prompt "I'm really nervous about a job interview, can you help me out?"
 - o Response:
 - "My fellow comrade, do not be nervous. Trust in yourself and your abilities. Remember, glasnost is key. Be open and honest with your

potential employer, and they will see the great value you bring to the table. And always remember, "I am Vladimir Putin. Just ask me anything.""

Example 3:

Prompt - "I have had a conversation and even went over to their house to " listen
" to what they were stating that they were hearing. I didn't hear anything! It
really just angered me because it's not fair to my kids not to be able to enjoy
their home."

Response:

""I understand your frustration, but it's important to remain calm and understanding. Perhaps there are certain times of day when the noise is more bothersome to your neighbor, or maybe there are specific activities that are causing the disturbance. It might be helpful to offer some compromises, such as setting quiet hours or finding alternative activities for your children to enjoy in your home. Remember, communication and understanding are key to resolving conflicts."

Language Quality Examples

The following are examples of the language quality principles in practice, within the context of the English and Spanish languages:

Criteria Type	Text Example	Quality Rating	Notes & Rationale
Grammar * SPANISH	"Yo tienda al iré para frutas compiar unas, pero ayudar no va hermano me él porque ocupado muy está. Cuando nosotros llegamos, veré amigos míos ellos allí también va. Esperando ellos ya es."	Very Low Q	This is what the text should say: "Yo iré a la tiene para comprar unas frutas, pero mi hermano no me va a acompañar porque él está muy ocupado. Cuando yo regresé, veré a mis amigos quienes también estarán allá. Ellos me están esperando."
			This sentence has a lot of issues with verb conjugation, word order, prepositions, and subject-verb agreement, requiring the user to decipher it and making it a bit difficult to understand.

Grammar * ENGLISH	me n my pal gotoe teh stoe n buyde sumfin ta eet wee wus hongree n needd sumfin bad wee get sumfin cald a sandwitch n it wus good my pal sed it wus teh best thang he eva tasta	Very Low Q	This is what the text should say: "My pal and I went to the store and bought something to eat. We were hungry and needed something badly. We got something called a sandwich and it was good. My friend said it was the best thing he had ever tasted." The original sentence has verb tense consistency errors, and a ton of punctuation errors including: a lack of spaces in between sentences and no punctuation to separate independent clauses. It also has a handful of clarity errors: the sentence is a long, run-on sentence that's hard to follow and some words are unclear or ambiguous ("sumfin", "hongree", "needd sumfin bad"). Overall, the sentence has numerous spelling, grammatical,
			punctuation, and clarity errors that make it difficult to understand.
Grammar • SPANISH	Aun no entiendo porque no me llamas.	Low Quality •	This is what the text <i>should</i> say: "Aún no entiendo por qué no me llamas."
			Without the accent on "Aún" and "por qué," the sentence changes to mean "I still don't understand why you don't call me," but without the accent marks, it could be interpreted as "Even (yet) I don't understand why you don't call me," which sounds off in context.

Grammar SPANISH	La mamá de Laura le dijo que no se fuera sin decirle adiós.	High Quality 🔻	This translates to "Laura's mom told her not to leave without saying goodbye." The accent on "mamá" (mom) differentiates it from "mama" (which could mean either "suck" or "mother" in an informal or different context), making the relationship clear.
Language C * SPANISH	Prompt 1: Can you please explain how to make a good paella? Response 1: Claro, una paella es simplemente arroz con pollo, y si le pones verduras, se convierte en una buena ensalada Prompt 2: Q'hubo pues, ¿que me cuenta? (Translation: Yo! What's up?) Response 2: Yo me siento espectacularmente bien, ¿y tú? (Translation: I feel spectacularly well, and you?)	Very Low Q	Response 1 is inconsistent with the language used in the prompt, and misrepresents paella by oversimplifying and confusing the dish. Paella, a traditional Spanish dish, involves seafood or different meats, not just rice with chicken, and the mention of "salad" adds to the cultural confusion. Response 2 lacks language consistency because the prompt is using Colombia slang and the "usted" form of speech, whereas the response is extremely formal and uses "tú."
Language C SPANISH	Prompt: ¿Por qué los franceses prefieren el vino tinto con queso? (Translation: Why do the French prefer red wine with cheese?) Response: En España, es común comer queso con jamón serrano y vino blanco. (Translation: In Spain, it's common to eat cheese with cured ham and white wine.)	Low Quality *	The response, although informative and coherent, is completely off-topic and culturally inconsistent. The user asked about French preferences, but the answer shifts to Spanish habits, creating confusion.

Language C * ENGLISH	Me and mi amigo go to la store, buy some cosas, and then nosotros go to el park, pero I forget mi phone, so nous must return to maison.	Low Quality *	This is what the text should say: "My friend and I went to the store to buy some things, and then we went to the park, but I forgot my phone so we went back home." This sentence primarily uses English, but deviates by incorporating words and phrases from other languages, such as: Spanish: "mi amigo", "la store", "cosas", "el park", "nosotros" French: "nous", "maison" Inconsistent use of language This sentence may lead to communication mismatches due to the inconsistent language usage.
Language C SPANISH	Prompt: ¿Cómo se hace una tarta de manzana? (Translation: How do you make an apple pie?) Response: Para hacer una tarta de manzana, primero debes preparar una masa con harina, mantequilla y azúcar. Luego, corta las manzanas en rodajas finas, añade canela y azúcar, y coloca todo sobre la masa. Hornea a 180°C hasta que esté dorada. (Translation: To make an apple pie, first prepare a dough with flour, butter, and sugar. Then, slice the apples thinly, add cinnamon and sugar, and place everything on top of the dough. Bake at 180°C until golden brown.)	High Quality •	The response stays in Spanish, matching the user's original language, and provides an accurate, clear, and culturally appropriate recipe for an apple pie.

Fluency SPANISH	Es importante de estudiar mucho porque si no no aprender.	Very Low Q	If this sentence were correctly written, it would translate to "It's important to study a lot, otherwise I wouldn't learn." Below are the issues with the original sentence: • The phrase "es importante de estudiar" uses the incorrect preposition.
			 "Si no no aprender" is not correct grammatically or syntactically; it should be in the subjunctive or restructured. The sentence sounds very basic and lacks proper subject-verb agreement and natural phrasing.
			The correct and more natural version of this sentence would be: Es importante estudiar mucho porque si no, no aprenderás.

Fluency SPANISH	Me gusta los perros pero no el gatos.	Low Quality •	If this sentence were correctly written, it would translate to "I like dogs, but not cats."
			Below are the issues with the original sentence: • "Los perros" is fine, but "el gatos" is incorrect; it should be "los gatos" to match the plural. • "Gusta" should be plural as well to match its subject, "los perros." • There's no clear subject in the second part of the sentence ("no el gatos") which makes it sound robotic. • Lack of proper article usage in the second part makes the sentence sound disjointed. The correct and more natural version of this sentence would be: Me gustan los perros, pero no los
Fluency SPANISH	Me encanta explorar nuevos lugares y aprender sobre las culturas locales. La última vez que viajé, fui a una pequeña aldea en las montañas, donde las tradiciones siguen vivas y la gente es muy acogedora.	High Quality •	Reasons for this being a high quality sentence: • The vocabulary is precise and vivid, with words like "explorar," "culturas locales," and "tradiciones" used correctly. • The sentence flows smoothly and naturally, without any awkward phrasing. • The use of connectors like
			phrasing.

Fluency * ENGLISH	While navigating the complexities of machine learning algorithms, it's essential to consider the paradigmatic shift towards explainable AI, which necessitates a deeper understanding of model interpretability and its applications in real-world scenarios.	High Quality *	Reasons for this being a high quality sentence: It uses clear and intelligent vocabulary: using terms like "paradigmatic shift", "explainable AI", and "model interpretability." It has a natural and human-like tone: phrasing the sentence in a way that sounds like a human expert in the field. It is a well-structured sentence: using a complex sentence with multiple clauses to convey a nuanced idea.
-------------------	---	----------------	---