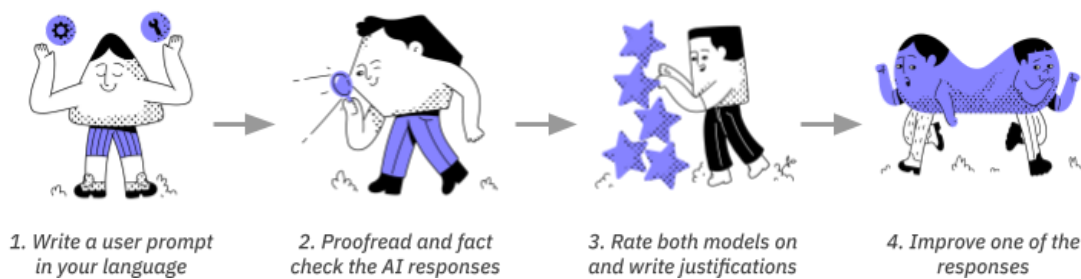


# Winter Wonderland RLHF | Do's and Don'ts

## ☰ Instructions | Winter Wonderland RLHF

### ✓ High-Level Tips

#### Task Overview



	Do!	Don't
<b>Prompt</b>	<ul style="list-style-type: none"><li>• Want to stick to the assigned category for the prompt</li><li>• Want the prompt to be specific and answerable</li><li>• Adhere to the desired complexity level (ie # of constraints)</li></ul>	<ul style="list-style-type: none"><li>• Choose a different category, you must stick to your category</li><li>• Have proofreading errors or poor writing</li><li>• Don't pick something out of your expertise or that is too hard to assess in terms of the response</li></ul>
<b>Response Ratings</b>	<ul style="list-style-type: none"><li>• Ensure that your dimensional ratings are objective and at least 9/10 of people would agree with you</li><li>• Fact-check every claim the model makes. Models often make up information!</li><li>• Pay attention to language fluency! The model will often use incorrect</li></ul>	<ul style="list-style-type: none"><li>• Not fact check!!!</li><li>• Assume the model is fluent in your language</li><li>• Don't select responses where theres a ton of complex/hard fact checking involved</li><li>• "Make up an error" → be intellectually honest in the way you rate the model</li></ul>

	<p>grammar, awkward phrasing or literal translations!</p> <ul style="list-style-type: none"> <li>• Spend time and be thorough on the ratings → make sure that you write your justifications in a way that helps the reviewer understand the issues you saw</li> <li>• Provide evidence for your ratings in the dimension-specific justification boxes</li> </ul>	
<b>Rewrite</b>	<ul style="list-style-type: none"> <li>• Remove pleasantries</li> <li>• Make sure the tone of your response is warm</li> <li>• Double check any factual information you are adding to the response</li> </ul>	<ul style="list-style-type: none"> <li>• Leave issues that you identified in your ratings unresolved (fix everything you identified)!</li> </ul>
<b>Side-by-side Rating</b>	<ul style="list-style-type: none"> <li>• Write a thorough justification explaining: <ul style="list-style-type: none"> <li>◦ Your verdict</li> <li>◦ Why the response is better</li> <li>◦ Specifically what it does better on</li> <li>◦ When you write your justification make sure to explain your thinking in detail (even if you have already explained it in the dimension-specific justification)</li> <li>◦ Your final justification should include everything you have mentioned in the dimension-specific justifications</li> </ul> </li> <li>• TRIPLE CHECK YOU'RE REFERENCING THE RIGHT RESPONSE (CHECK THAT RESPONSE 1 AND RESPONSE 2) ARE WHAT THEY ARE.</li> <li>• Make sure your side-by-side rating</li> </ul>	<ul style="list-style-type: none"> <li>• Be generic in your justification</li> <li>• Make your side-by-side rating inconsistent with your dimension justifications (ie if one response is rated higher on most dimensions you should not prefer the other response on your side-by-side rating)</li> </ul>

	aligns with your dimensional ratings (ie if one response is rated higher on most dimensions, it should generally be preferred on the side-by-side rating)	
--	---	--