

Assessment Method Used	Characteristics of Incorrectly Answered Questions
BERTscore	<ul style="list-style-type: none"> <li>• Listing Elements</li> <li>• Repeating Incorrect Elements</li> </ul>
sacreBLEU	<ul style="list-style-type: none"> <li>• Summarization Tasks</li> </ul>
Manual Error Analysis	<ul style="list-style-type: none"> <li>• Answering multiple questions</li> <li>• Mismatch of Nomenclature</li> <li>• Inferring Numerical Operations from Attribute Names</li> <li>• Ambiguous Questions</li> </ul>