

Exploring the Development of Briefing Books for Early Scientific Advice Using Large Language Models: A Proof-of-Concept Study

Ryan Thalifffdeen¹, Matthew Radford¹, Inês Guerra², João Leite³, Shaktidhar Pullagurla⁴, Vishalie Shah², Edel Falla², Raja Shankar², Yumi Asukai⁵

¹Gilead Sciences, Inc., Foster City, CA, USA; ²IQVIA, London, UK; ³IQVIA, Lisbon, Portugal; ⁴IQVIA, Bangalore, India; ⁵Gilead Sciences, Ltd., Stockley Park, UK

Copies of this poster obtained through QR (Quick Response) and/or text key codes are for personal use only and may not be reproduced without written permission of the authors.



Conclusions

- The LLM performed better in areas that relied on pre-trained knowledge (e.g., choice of comparator), compared to areas that required advanced reasoning (e.g., ITC, economic modeling)
- Overall, the LLM, although successful in retrieving information from the knowledgebase, could not generate an HTA-grade BB
- Improving the knowledgebase with relevant literature and clinical feedback, coupled with expert prompting guidance, could enhance the LLM's performance. However, this improvement would come at the cost of considerable human effort

Introduction

- Health Technology Assessment (HTA) bodies (e.g., National Institute for Health and Care Excellence (NICE)) are market-specific groups that evaluate the clinical, safety, and economic evidence surrounding a new medicinal product coming to a local market¹
- Early Scientific Advice (ESA) is an opportunity to inform Clinical Development and Reimbursement Strategies.² Prior to any ESA engagement, a briefing book (BB) will be developed to seek advice on particular topics and summarize the company's position on those questions
- While Large Language Models (LLMs) have demonstrated efficiencies for various Health Economics and Outcomes Research deliverables, BBs pose unique challenges for LLMs, as BBs are generated earlier in a product's lifecycle when evidence is limited
- Additionally, BBs require strategic thinking to develop a company's position and justification on questions for HTA in order to optimize the BB to the specific necessities of each HTA process

Objective

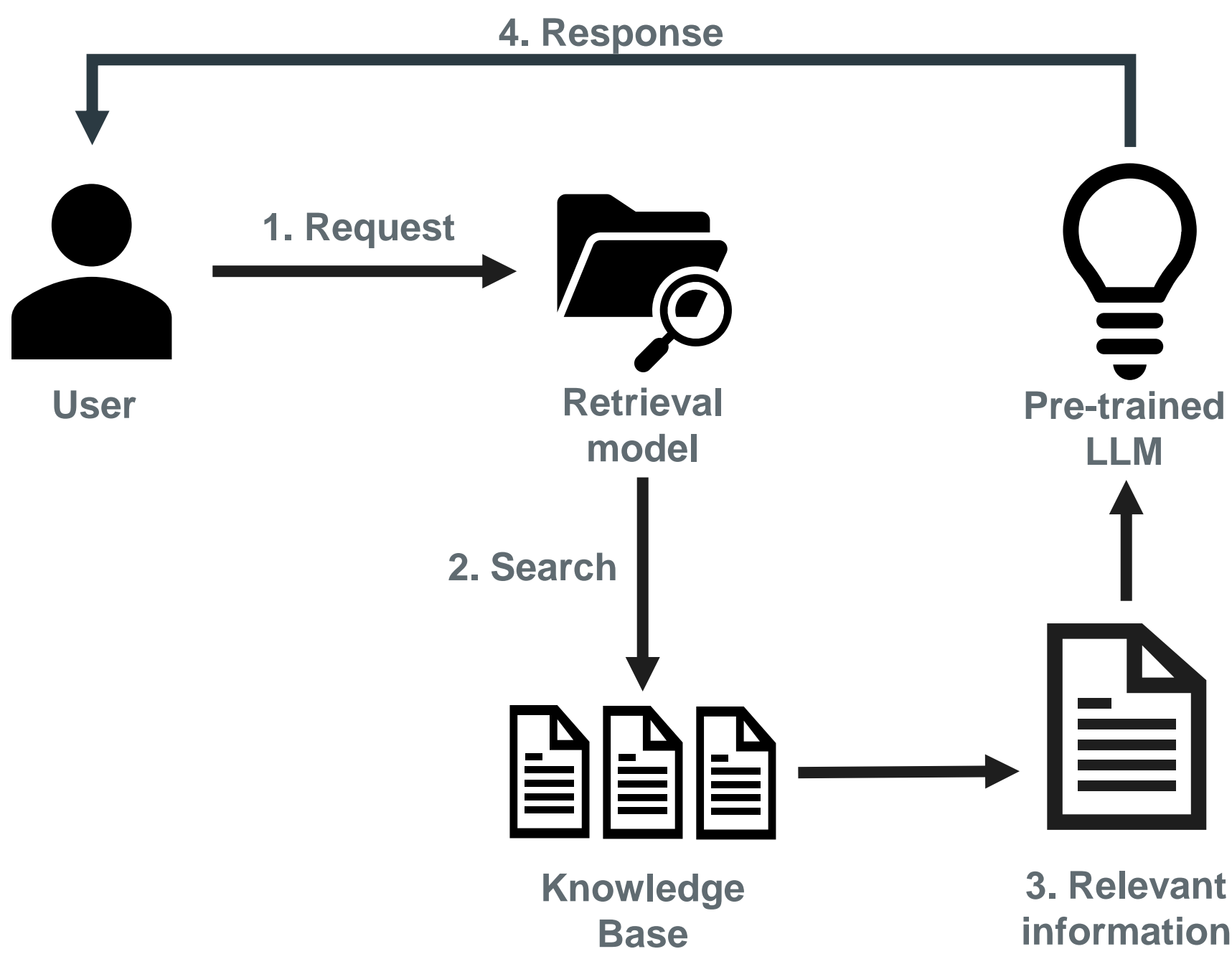
- This proof-of-concept study aimed to assess the feasibility of LLM-based generation of BBs for ESA

Methods

Sections and Knowledge Base

- Sections on approaches for trial comparator selection, indirect treatment comparison (ITC), and economic modeling were created using GPT-4 via python API
- To supplement the model's pre-trained knowledge, retrieval-augmented generation (RAG) was used for content generation and answer retrieval (**Figure 1**)
- The model's knowledgebase included the trial protocol, internal strategic documents, previous HTA appraisals, HTA BB guidance, and published trial results in similar indications
- Prompts were developed iteratively upon review of outputs
- Key Evaluation Metrics were output quality and human-led effort needed for revisions

Figure 1. Retrieval Augmented Generations (RAG) Framework



Results

Figure 2. The first question focused on the justification for the choice of comparator in the trial

Question	
Does the Agency agree that neoadjuvant nivolumab in combination with chemotherapy, followed by investigator's choice adjuvant treatment, is an appropriate comparator representing the SOC for patients with early-stage solid tumour?	
Prompt	Results
Prompt with no sources was tested as control	LLM retrieved the message "The information needed to answer the question is not provided in the context. "
Initial prompt introduced a clinical trial of nivolumab as a source and requested the model to utilize NCCN guidelines for justification	Model provided short summary of CheckMate 816 trial and hallucinated using the NCCN guideline as justification
Request to use NCCN guideline was removed from the prompt	Model provided more extensive summary of CheckMate 816 trial and limited top-level information regarding standard of care.
Revised prompt included two main updates : (1) additional knowledge sources (HTA recommendation for Nivolumab, HTA recommendation for Atezolizumab, HTA briefing book template and briefing book guidance); and (2) a clearer distinction between the neoadjuvant and adjuvant setting within the context.	The response from the LLM was generally similar to earlier responses , although the LLM was able to identify that the investigator's choice of therapy is consistent with country-specific guidelines in the adjuvant setting. While the model was able to recognize the relevant sources , it did not always retrieve all of the relevant information .

Results

Figure 3. The answers by the LLM were often high level and did not always utilize the relevant information from the given sources

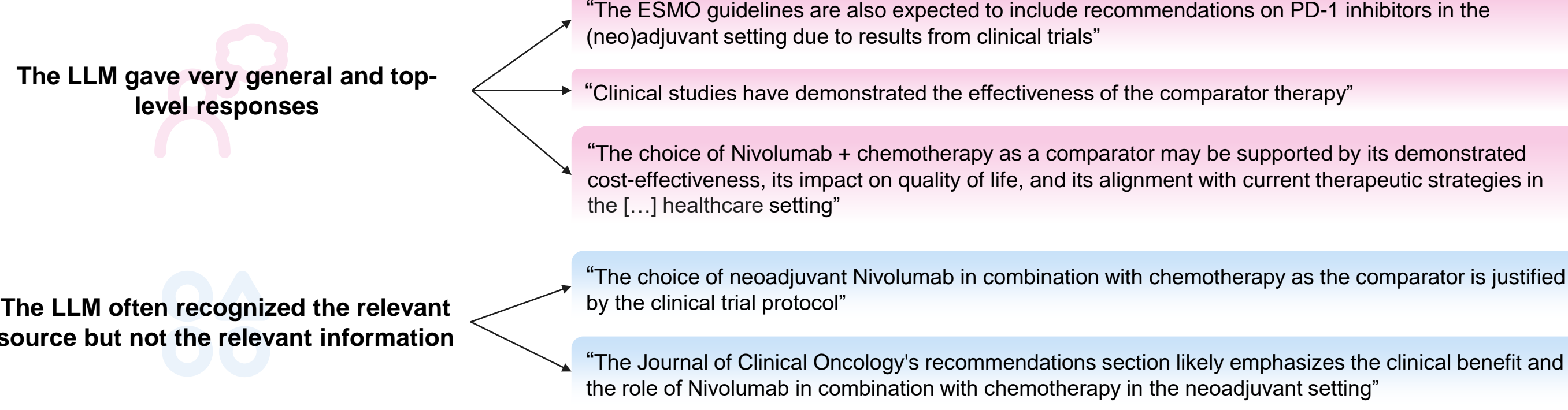


Figure 4. The second question focused on whether the LLM could generate and justify an economic model based on the relevant context provided

Question	
Does the Agency agree that the proposed economic model structure captures the key health state relevant for this patient population, and do the available trial data/endpoints allow those health states to be accurately populated over time? Is there an alternative model structure/approach that should be considered?	
Prompt	Results
Initial prompts for the RAG model included providing a summary of the economic model and requesting for an expansion on the model design and any considerations from the LLM	The answer by the LLM regurgitated the summary of the economic model in different words, providing no new information
Another approach explored using pre-trained knowledge to provide context on the objective of the model and request the LLM to provide a suggested model design for the objective and clinical context.	The LLM suggested slightly different health states in different attempts , and it proved inconsistent in terms of model design. Some of the responses were more aligned with the health states provided than others. Additionally, the model was not capable of conceptualizing specific design aspects for the decision problem
A tutorial on economic models was provided to the LLM in order to provide further context and options for model designs	The tutorial did not improve the response , and introduced some confusion to the model , which took the model examples from the tutorial and utilized their approach for the model design

Figure 5. As the response by the LLM was not satisfactory, two different approaches were explored to further elaborate on the economic model approach and design suggested by the LLM

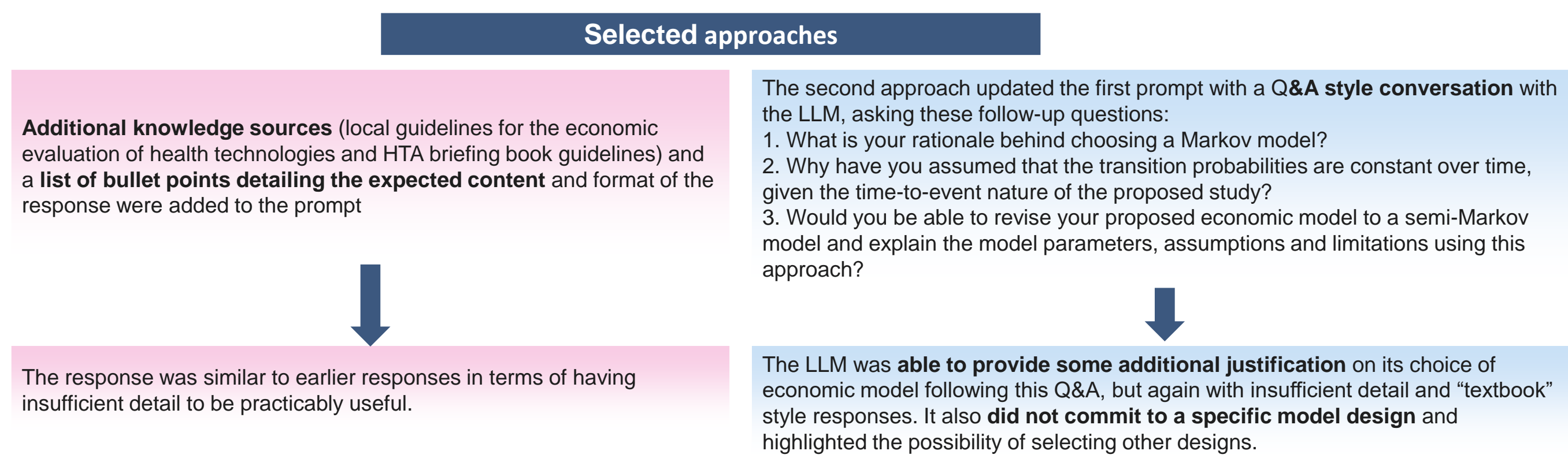


Figure 6. The LLM could not conceptualize an economic model beyond the basic assumptions and standard model designs

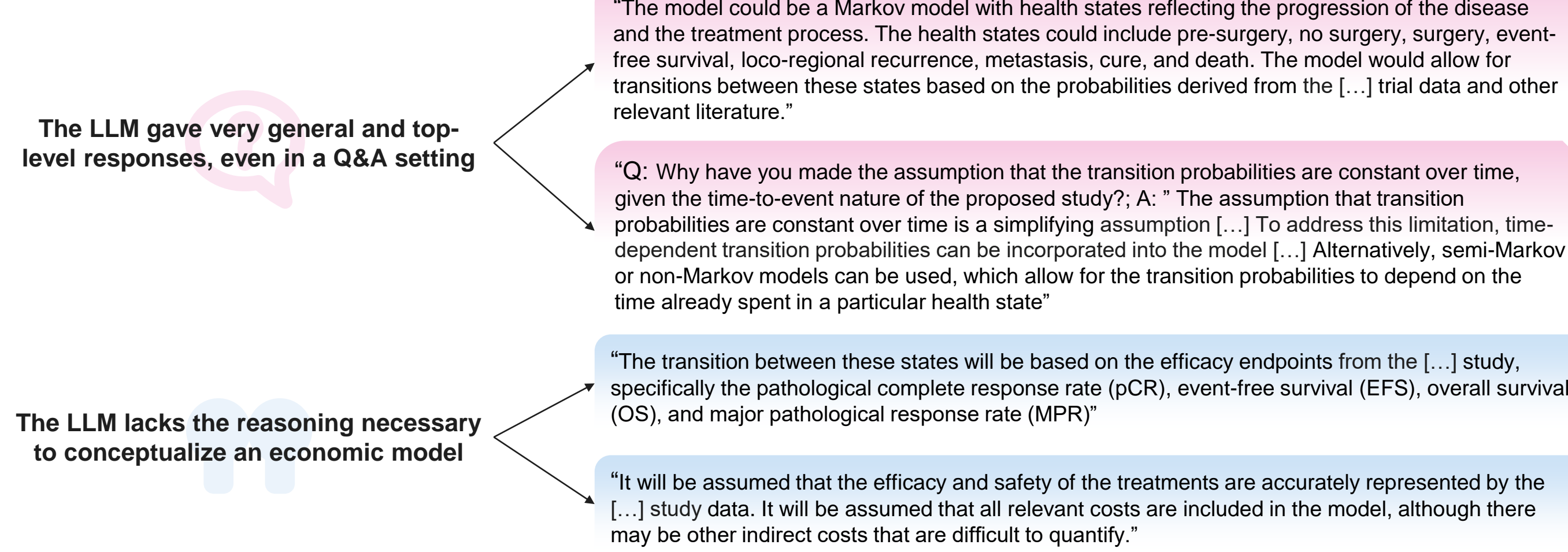
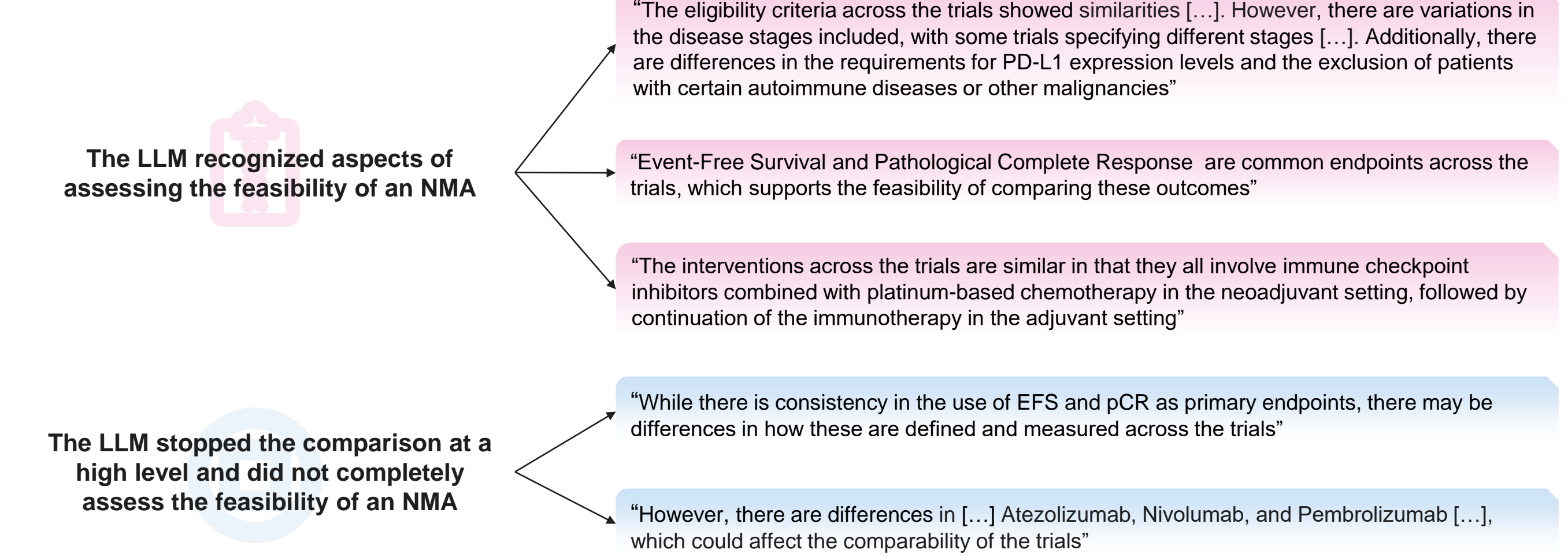


Figure 7. The third question focused on whether the LLM could assess the feasibility of conducting an NMA/ITC based on the pivotal trial design

Question	
Can the Agency comment on the feasibility and utility of conducting a network meta-analysis (NMA)/ITC to compare the EFS and pCR of product X with other relevant comparators?	
Prompt	Results
Initial prompts included only the Nivolumab HTA report and trial protocol as the sources	The answer by the LLM was top level and non-specific , and it did not provide any assessment on the possibility of an NMA
One solution tested was to provide the LLM with a list of possible relevant comparator trials and retrieve their relevant information for comparison with the pivotal trial via clinicaltrials.gov. Only Inclusion/Exclusion Criteria, Arms and Interventions and Endpoints were compared.	Response was more positive, and it provided a general comparison between the trials , their characteristics and a short assessment on the feasibility of conducting an ITC. However, the LLM failed to recognize most of the necessary steps and criteria for conducting an NMA
Two clinical trial publications were given to the LLM in order to assess a possible network with the trial, instead of the Clinical Trials website. This ideally allowed the LLM to assess the reported baseline characteristics that were not reported in the website.	Response was similar to previous ones. The LLM provided a short analysis on the feasibility of comparing study populations, possible subgroups of interest, differences between trials, and a possible anchor comparator for an ITC. The LLM lacked the reasoning capabilities to provide any justification for its assessment and it did not compare baseline characteristics.

Figure 8. The LLM recognized some criteria necessary for an NMA but lacked the depth to properly assess the feasibility of an ITC



References: 1. Trowman R et al. Health technology assessment 2025 and beyond: lifecycle approaches to promote engagement and efficiency in health technology assessment. International Journal of Technology Assessment in Health Care. 2023;39(1):e15. doi:10.1017/S0266462323000090. 2. Wang T, et al. Building HTA insights into the drug development plan: Current approaches to seeking early scientific advice from HTA agencies. Drug Discov Today. 2022 Jan;27(1):347-353. doi: 10.1016/j.drudis.2021.09.014. Epub 2021 Sep 28. PMID: 34597755.

Correspondence: Ryan Thalifffdeen (Ryan.Thalifffdeen@gilead.com)