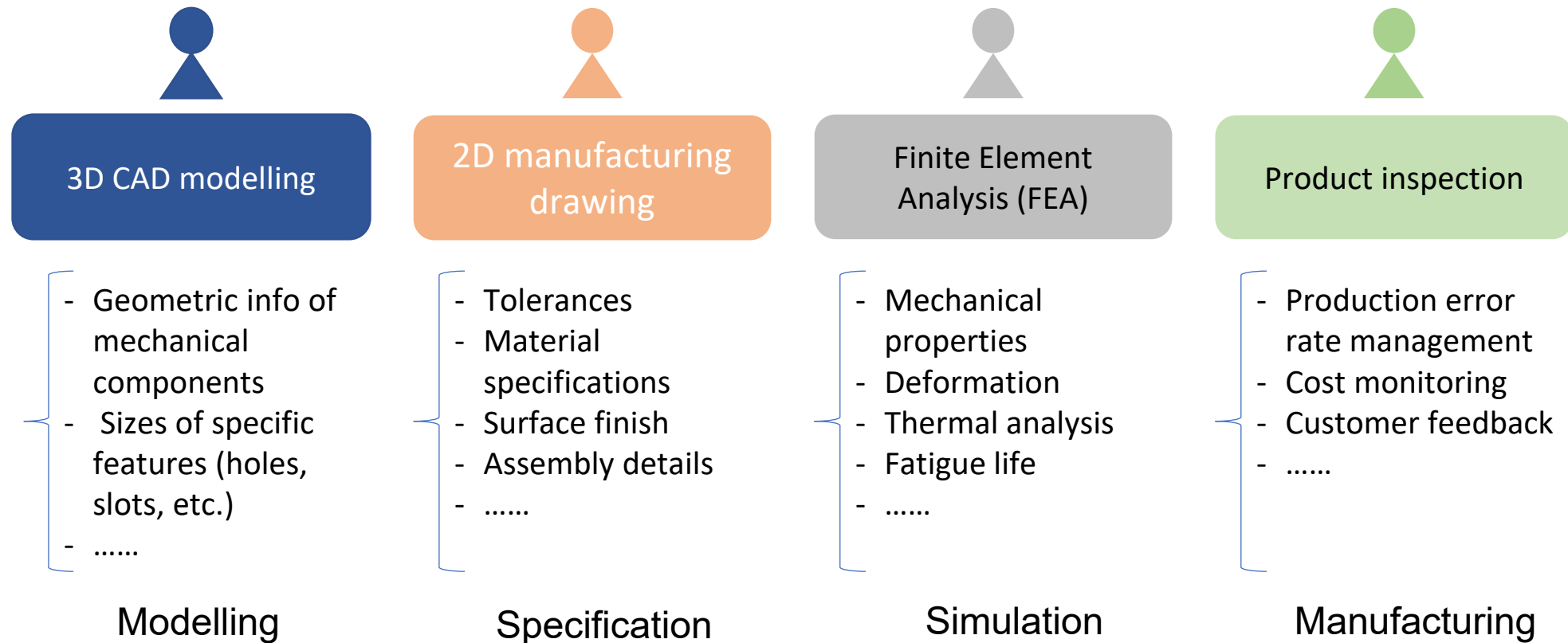


MechRAG: A multimodal large language model for Mechanical Engineering

Shuang Li, Jonathan Corney
School of Engineering, University of Edinburgh
09/2025

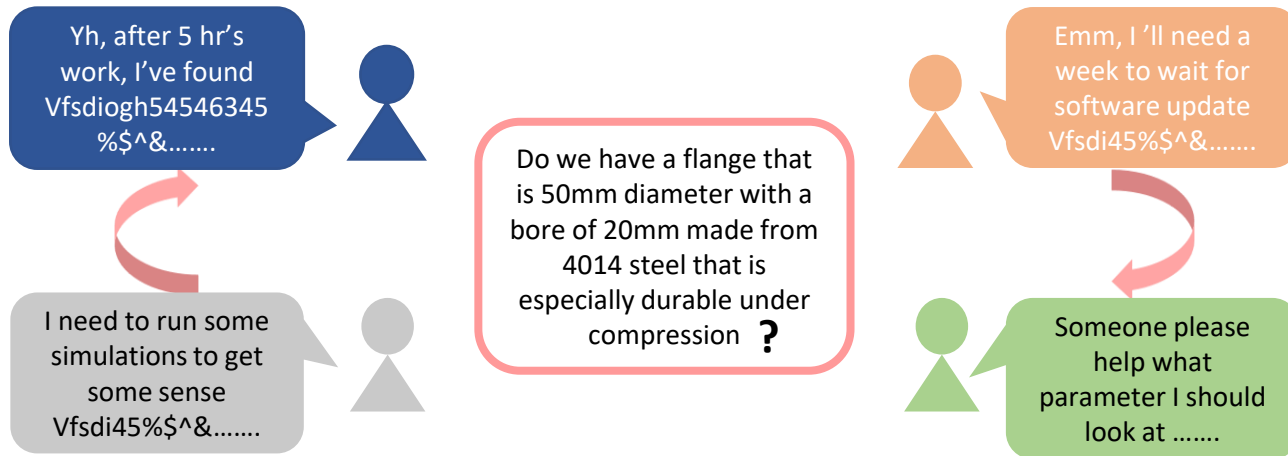


Challenge in Technical Design: CAE Software forms distributed silos of data and knowledge

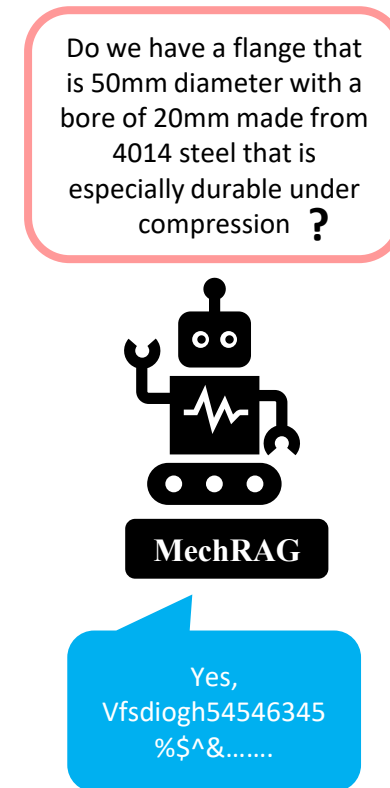


Challenge in Technical Design: CAE Software forms distributed silos of information and knowledge

Current Practice: The answering to technical queries frequently requires information from multiple silos.
Current solutions are facilitated by forms of team conversations.



Research vision of this work: AI responds to queries using holistic knowledge from multiple silos



Classifying Technical Design Questions

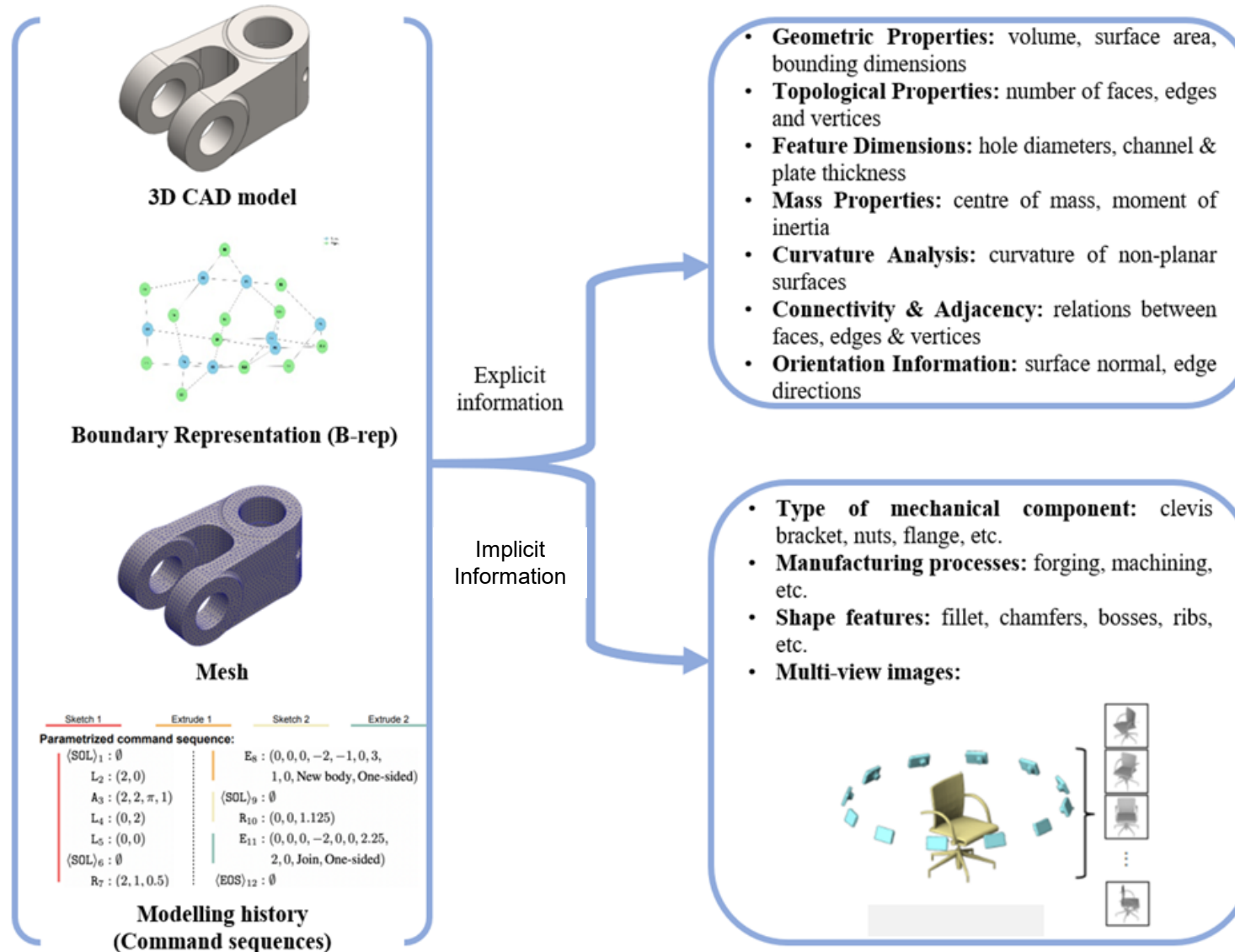
Level	Type	Illustrative Examples of MechRAG Prompts	
1	Direct	<p>Query information of single parts</p> <ul style="list-style-type: none"> What is the volume of part 00210008 in cm³? How many round holes does model 00210337 contain? Does part 00210217 have any type of symmetry? 	Multi-silo queries
2	Collective	<p>Query combines Level 1 information from multiple parts</p> <ul style="list-style-type: none"> Which parts have the biggest volume and the smallest? What symmetries are present in the dataset of components? 	Multi-silo, multi-design, comparative queries
3	Emergent	<p>Query requires combined information from both CAD data, and the general knowledge of LLMs obtained from pre-training</p> <ul style="list-style-type: none"> Based on part 00217697's modelling history which other components is most similar to it? Which of 0021* series of components would be the most expensive to produce and why? What manufacturing processes is used to produce part 00210215? 	Subjective queries beyond given data

Tokenising CAE Data

CAE data is typically numerical, so it needs to be processed to provide LLM with compatible formats (e.g. images and text).

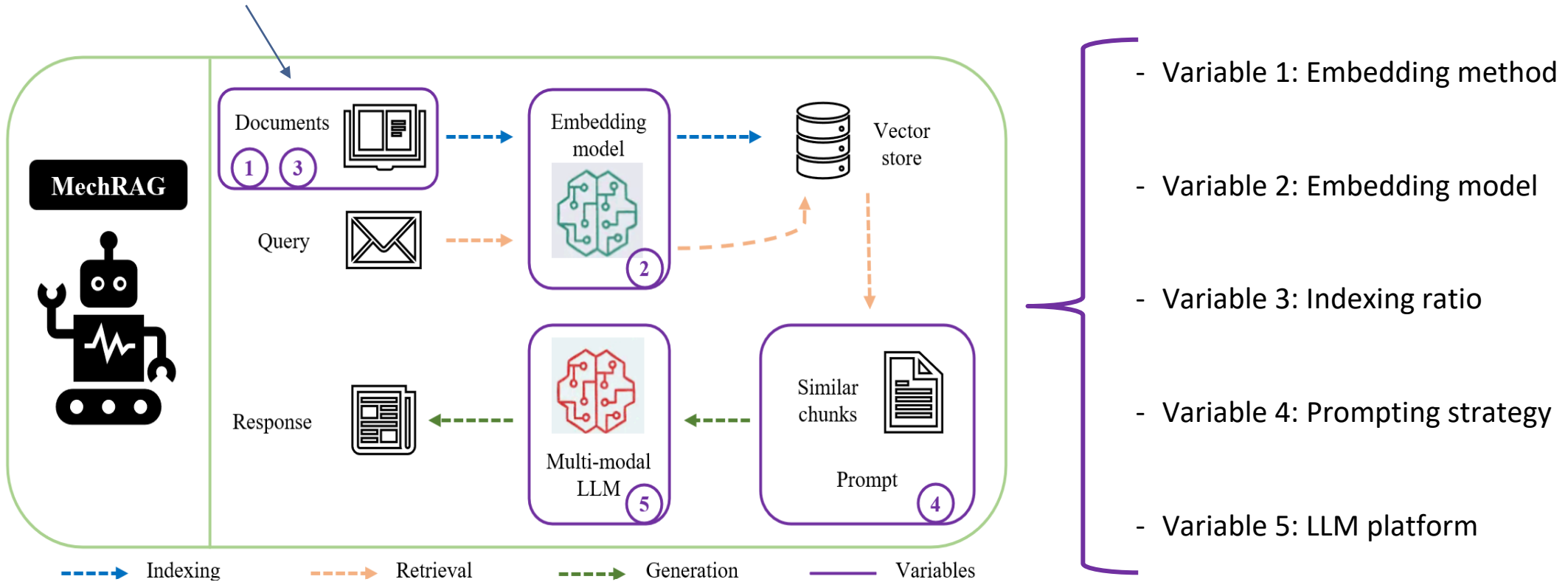
The volume of the model_00210231 is 237671.829 cubic MM.
 The surface area of the model_00210231 is 38178.907 square MM.
 The volume of this model_00210231 bounding box is 374883.431 cubic MM.
 The variance of the Gaussian curvature of this model_00210231 is 0.008.
 The variance of the Mean curvature of this model_00210231 is 0.000.
 The flat plane of this model_00210231 occupies 0.287%.
 The curved surface of this model_00210231 occupies 99.713%.
 The straight edge of this model_00210231 occupies 24.490%.
 The curved edge of this model_00210231 occupies 75.510%.
 This model_00210231 is radially symmetrical.
 The axis of symmetry of this model_00210231 is [-2.98027006e-10 -1.00000000e
 One of the round holes in this model_00210231 has a diameter 1 of 60.000 MM
 One of the round holes in this model_00210231 has a diameter 2 of 9.000 MM
 One of the round holes in this model_00210231 has a diameter 3 of 9.000 MM
 One of the round holes in this model_00210231 has a diameter 4 of 9.000 MM
 One of the round holes in this model_00210231 has a diameter 5 of 9.000 MM
 One of the round holes in this model_00210231 has a diameter 6 of 60.000 MM

Fragment of a **text summary** that combines information from different CAE representations (e.g. BRep, Mesh, Design History etc).



MechRAG architecture

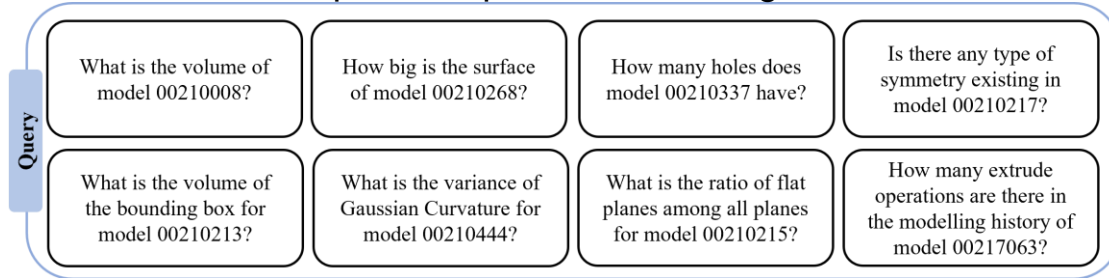
Multi-modal summary
of CAE data



The **Retrieval-Augmented Generation (RAG)** mechanism enhances large language models (LLMs) by combining their generation capabilities with external information retrieval. Instead of relying solely on their internal training data, RAG models first retrieve relevant documents or knowledge snippets from an external database or knowledge base based on the input query. They then use this retrieved information to ground and inform their generated responses, leading to more accurate, up-to-date, and contextually relevant outputs without the need for retraining the LLM itself.

MechRAG example generations

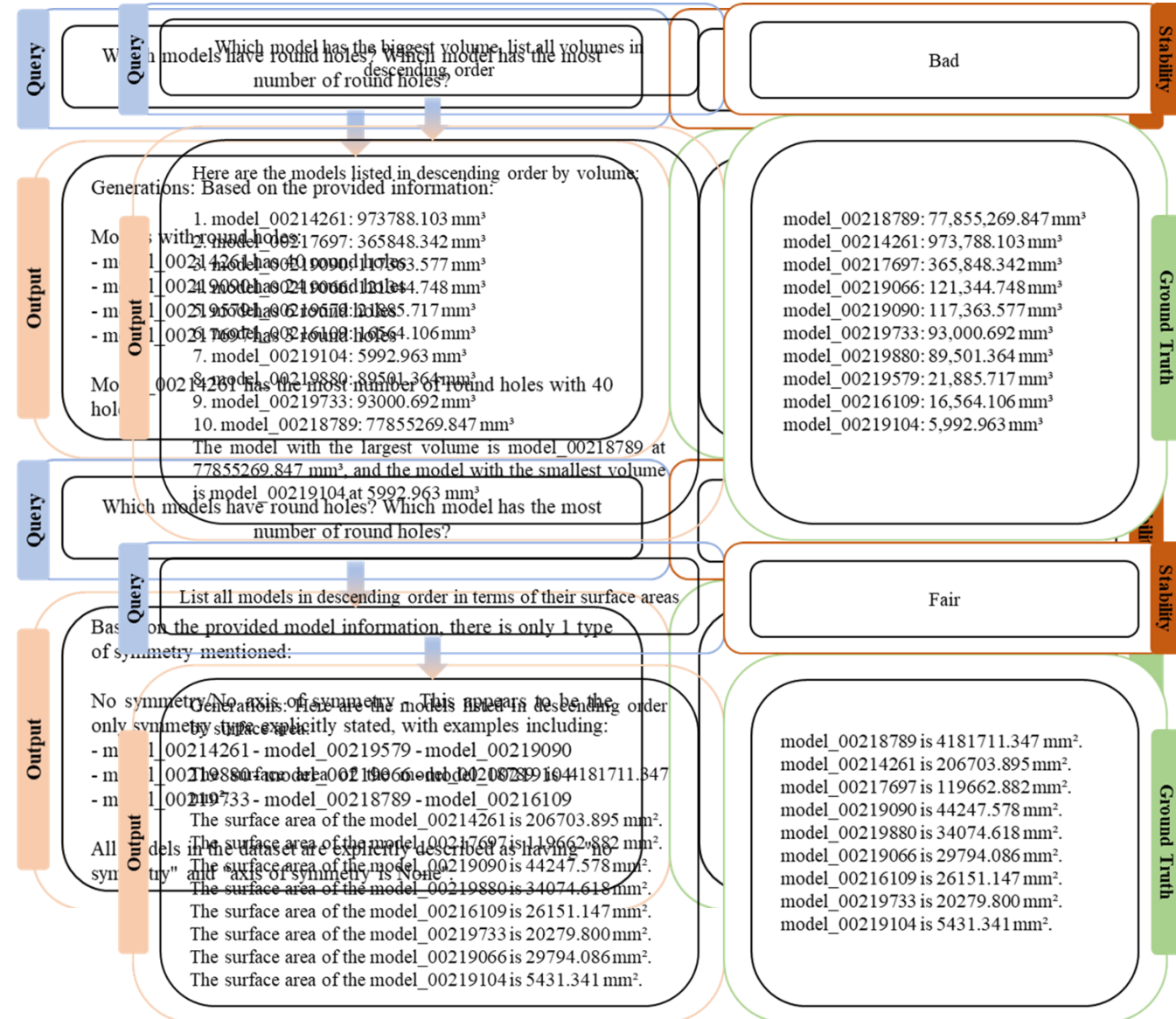
Level 1 queries: specific info of single models



For Level 1 queries MechRAG achieved an accuracy of 90.48%.

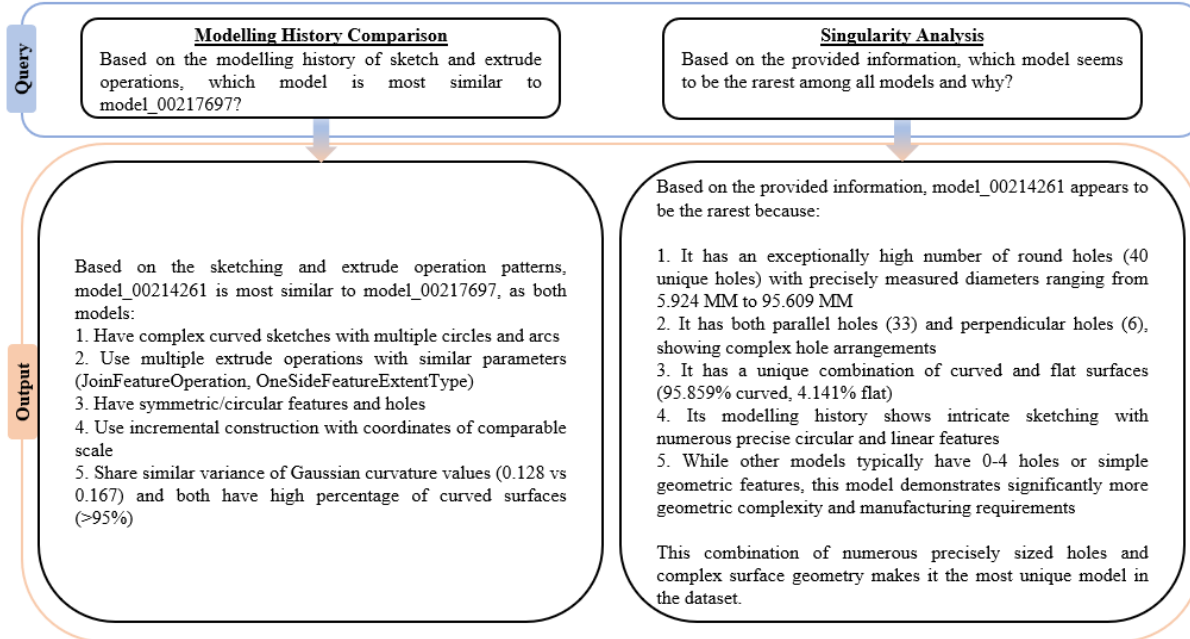
For Level 2 queries, performances decline when the quantity of involved numeric values increases.

Level 2 queries: analysis of groups of models

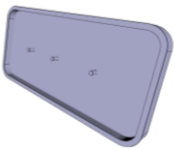


MechRAG example generations

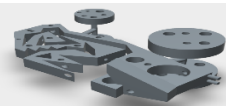
Level 3 queries: ambiguous prompts “most similar” & “rarest”



00217697



00214261

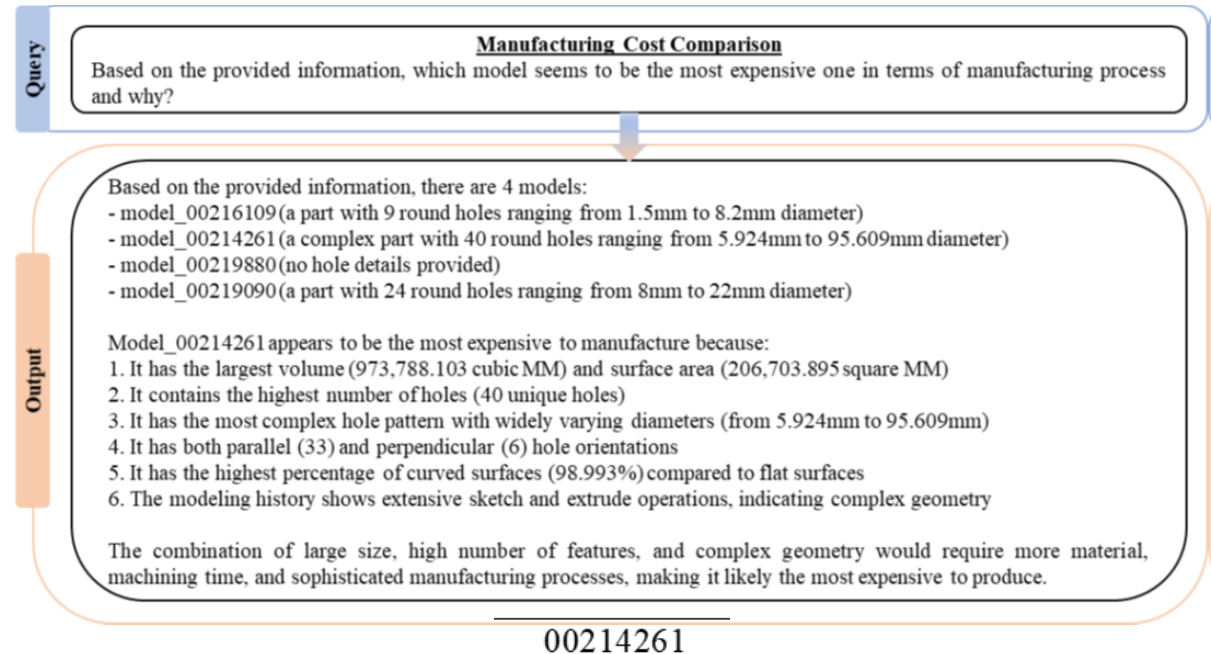


00214261



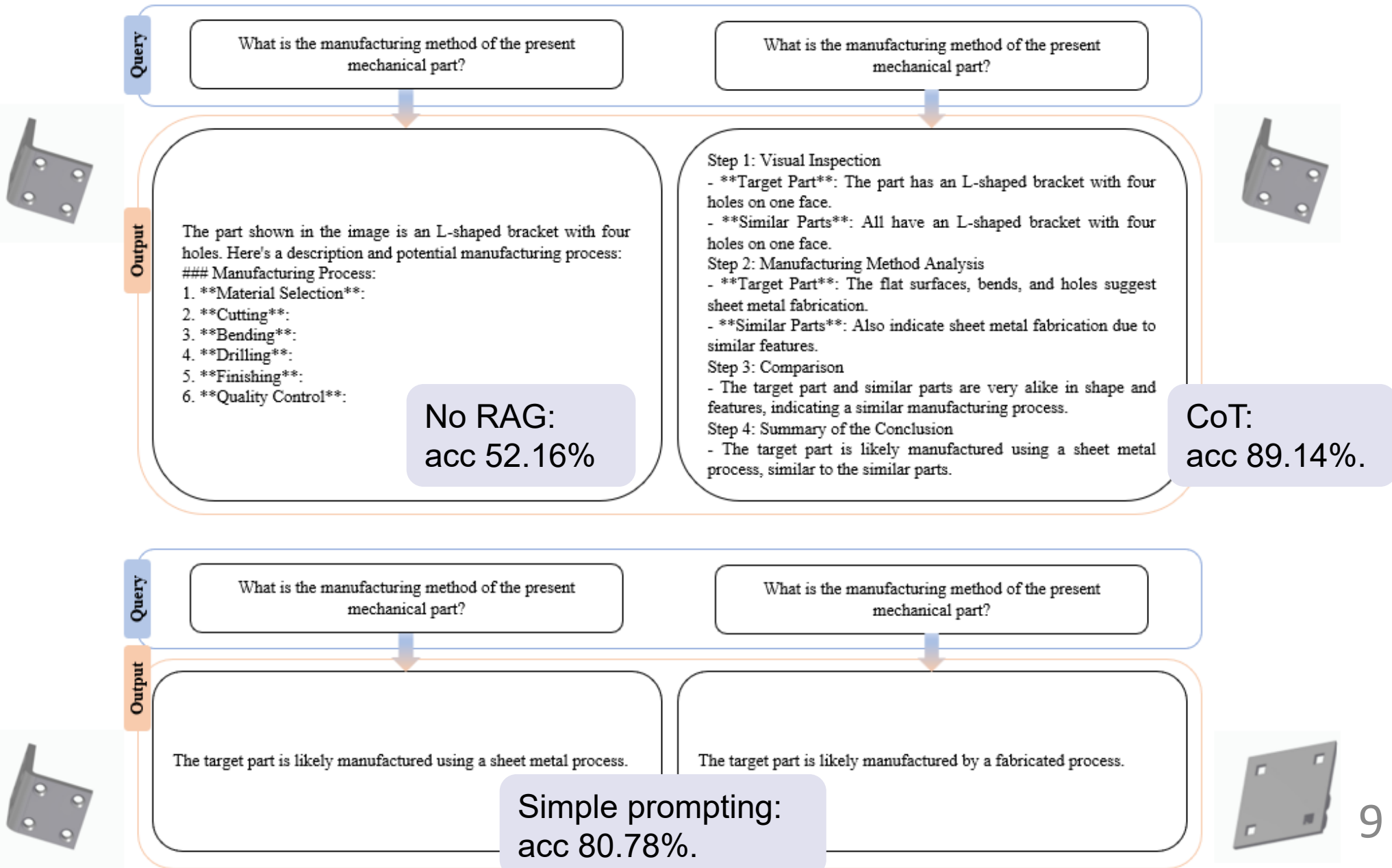
For Level 3 queries, MechRAG responds with mix of CAE data (via RAG) and general knowledge, and demonstrates satisfactory reasoning capabilities.

Level 3 queries: ambiguous prompts “most expensive to manufacture”



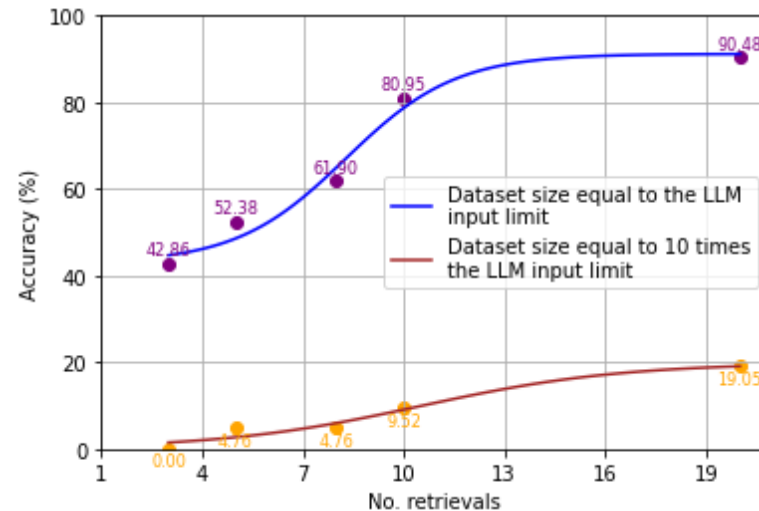
MechRAG example generations

MechRAG can also manage multi-modal conversations. The best accuracy achieved for manufacturing process recognition is 89.14%.



Future Work

- ❖ **Token limits of the context window in LLMs:** it's not unusual for engineering companies have tens of thousands of CAD models, but it's possible to exceed an LLMs token limit with less than a hundred models. Although the token limits of LLMs are expected to increase, future work is needed to minimise the number of tokens required to represent engineering product knowledge.



- ❖ **How to utilise other modalities/formats of data** (e.g. FEA/CFD simulation and process) to augment LLMs?

Help us! – Complete Our Quick Survey

- Link to the survey:

<https://forms.office.com/Pages/ResponsePage.aspx?id=sAafLmkWiUWHiRCgaTTcYanOEx6JrOI-FqdKn14ZKk9NUMlpHODQ0Vk5VNFFQQINEMkhINUM1QjBBNy4u>

Survey on MechRAG performances



Thanks for your attention!

