{"NCT_ID": {brief_title: "The title of the study ", phase: "phase the trial is in", drugs: str(list()) (single string), drugs list: ["drug1", "drug2", ..., "drugn"] {"_id": "sigir-20141", "text": "A 58-year-old diseases: str(list()) (single string) {"_id": "sigir-20142", "text": "An 8-year-old {"_id": "sigir-20143", "text": "A 58-year-old diseases_list: ["disease1", "disease2", ... , "diseasen"] _id": "sigir-20144", "text": "A 2-year-old enrollment: str(number of people enrolled) inclusion_criteria: "inclusion criteria having \n\n as a separator" exclusion_criteria: "exclusion criteria having \n\n as a separator" brief_summary: "summary of the trial if given" **STEP 1 – KEYWORD GENERATION** Using gpt-4-turbo, user prompt given to generate upto 32 keywords from the patient raw case study. Clinical trial query file (corpus.jsonl): json line format, processed trials information for efficient retrieval Query Patient file (id2quires.json): json format { id: "NCT ID", title: "brief title of the study", text: "brief summary + inclusion + exclusion", metadata: "all info present in raw file" } {"Patient id": {raw: "patient case study", model_name:{ Sample: **summary**: "patient summary by model used", **conditions**: ["C1", "C2",..., "Cn"] (list of keywords extracted from patient upto 32 keywords) RETRIVAL PROCESS (ALL TRIALS RELEVANT TO PATIENT) Input: The "text" field of the trials only. Search algorithm: Hybrid search Input: each keyword of the given patient. Trial query conversion to sparse vectors: BM25 ranking model - BM25Okapi, stored as tokenized json format. Query conversion to vectors: MedCPT-Query-Encoder model. Used for vectorization of smaller texts. Trial query conversion to dense vectors: Step 1 is to tokenize text using "AutoTokenizer" MedCPT-Article-Encoder model. Used for vectorization of large texts. - Step 2 is to convert the tokens into embeddings using "AutoModel" - Step 1 is to tokenize text using "AutoTokenizer" Step 2 is to convert the tokens into embeddings using "AutoModel" Ouput: a pytorch tensor of length 256. Ouput: a pytorch tensor of length 512. Storage: "Faiss" for MedCPT embeddings. Input parameters : **K** = int() "different k for fusion", bm25_wt = int() "bm25 weight (the weight used here was 1), medcpt_wt = int() "MedCPT weight (the weight used here was 1), N = int() "top k matches to be retrieved for each query and at final aggregation". BM25 search output: Top N list of NCT_ids for every key words for a patient based on ranking. MedCPT search output: Top N list of NCT ids for a every key words for a patient based on embedding match. **Calculation of BM25 and MedCPT scores:** Scoring function = (1/(rank+K))(1/(combination_index+1)) **Combination of BM25 and MedCPT scores:** The score for a NCT_ID is 0 at start. Considering a keyword match with NCT1 in BM25 scoring 0.5 and the same keyword is matched with NCT1 in MedCPT retrieval giving the score 0.2. Thus, a final score for NCT1 given the keyword hence the patient is 0.7 due to maximization problem. A json format, each patient having Top_N NCT_IDs based on scoring system in list: qid2nctids_results_{model_name}_{corpus}_{k}_{bm25_wt}_{medcpt_wt}_{N}.json {"P1": ["NCT1", "NCT2", ..., "NCT{N}"], "P2": ["NCT1", "NCT2",..., "NCT{N}"], ..., "Px": ["NCT1", "NCT2", ..., "NCT{N}"]} MATCHING PROCESS (PATIENT - CRITERION LEVEL SCORES) Output: The predicted reasonings for patient-criterion level as json format. sigir-20141": { **LEVERAGING LLM CAPABILITY (GPT-4):** "NCT01660594": } - customized prompt with details based on inclusion/exclusion criteria given Input: Processing of the output file from retrieval step as below Relevance of criterion contains output details and format output details: for each criterion, 1. Reasoning process of the criterion w.r.t. patient. 2. List of indices of relevant sentences w.r.t to criteria found in patient case. 3. Label for the criterion based on the reasoning, {"not applicable", "not enough "1": [information", "included", "not included"} "There is -output format: dict{str(criterion_number): list[str(element_1_brief_reasoning), [], list[int(element_2_sentence_id)], str(element_3_eligibility_label)]} "diseases_list": [**Processing of the patient data:** Case study converted to a sentence_index format using sentence tokenizer and List of indices having enrollment": "110.0", concatenated again to make whole of of patient study for input into the LLM. relevant sentences for "brief_summary": "To determine the "NCTID": "NCT01724996" the given criterion. User prompt: The patient details: "Processed sentence indexed form" Clinical trials: "trial criteria information based on the NCT_ID retrieved for the patient "exclusion": { using the file in left" Output: "how the output be represented, as a plain ison" Label for the criterion RANKING PROCESS (PATIENT - TRIAL LEVEL SCORES) **LLM AGGREGATION SCORE (GPT-4):** System prompt: - given the inputs, provide output of 2 scores. 1. Relevance Score (R) 2. Eligibility Score (E) output details: for each trial, 1. Provide an explanation for the relevance of the patient to the trial. 2. Provide a relevance score (R). R can be range from [0,100], - If R = 0 the patient it totally irrelevant to the given trial - If R = 100 the patient is exactly relevant to the given trial 3. Provide an explanation for the eligibility of the patient to the trial. 4. Provide a eligibility score (E). E can be range from [-R, R], - If E = -R the patient is ineligible (not included by any inclusion criteria, or excluded by all exclusion criteria) - If E = R the patient is is eligible (included by all inclusion criteria, and not excluded by any exclusion criteria) - If E = 0 the patient is neutral (i.e., no relevant information for all inclusion and exclusion criteria) -output format: Dict{"relevance_explanation": Str, "relevance_score_R": Float, "eligibility_explanation": Str, "eligibility_score_E": Float} **Processing of the patient data:** Case study converted to a sentence_index format using sentence tokenizer and concatenated again to make whole of of patient study for input into the LLM. **User prompt/inputs:** The patient details: "Processed sentence indexed form" Clinical description: "trial tittle along with disease conditions and a brief summary of the trial" Patient criteria predictions: "String format of the prediction json output generated in matching process" Output: "how the output be represented, as a plain json" **LINEAR AGGRIGATION SCORE:** FINAL OUTPUT: Matching score: Based on labels of the predicted criterion the scores is being updated via a Patient ID: sigir-20141 common match and penalty scoring system. Clinical trial ranking: NCT00185120 2.8999999999 Aggregation score: NCT02144636 2.8999999995 Aggregating the relevancy score and eligibility score to obtain a final aggregation score NCT02608255 2.84999999975 $aggregation_score = (R + E)/100$ NCT01724996 2.7999999998

Final ranking of trial for a given patient:

Rank = match_score + aggregation_score

Clinical trial raw file (trial info.json): json format, extracted and parsed through clinicaltrial.gov

Raw Patient file (quires.jsonl): json line format

{_id: "Patient id", text: " Containing patient case"}