

Question	<p>Answer</p> <p>Priya and Kunal, EXCELLENT answers overall :) My comments/corrections are in bold. #3 and #17 are the only places where my answers override yours :) You can DELETE this paragraph before sending this to the other CPs/TAs, and, feel free to remove my bolding of the text in my responses, as well.</p>
<p>1. When we upload pics/clips/songs to social media, what specific mechanism do we use to help others find our content when they search? Why do we need this specific mechanism?</p> <p>(1 Mark)</p>	<p>We use #hashtag(s) [one or more] so that anyone searching using a hashtag can find our content as well, eg. #MeToo.</p> <p>Hashtags are easy to index. They serve as textual descriptors for non-textual content even (eg songs), making search simple [because similarity search isn't widely implemented yet, it's brand new].</p> <p>Indexing of pictures, clips, songs and other media is usually done using metadata. (+0.5)</p> <p>Pictures, clips, videos, etc media do not exactly have textual data that will help the search engine index it. Hence the metadata, like description, can provide more information about the actual content. (+0.5 for similar explanation)</p>
<p>2. What algorithms causes/leads to/results in/is implicated in ... "filter bubbles"? How does it lead to this?</p> <p>(1 Mark)</p>	<p>"Filter bubbles" are caused by algorithms in personalized content recommendation systems. (+0.5)</p> <p>Recommendation systems suggest content that is tailored to the user's interests and preferences. Hence, the user is limited to seeing information that is restricted in the areas of their interests causing these "filter bubbles". (+0.5 for similar explanation)</p>
<p>3. What two other items can a search engine serve us (eg. via "snippets"), in addition to</p>	<p>[The question asks for new items... Any new items not currently being served will</p>

<p>what gets served already? Name each item, and briefly state why it would be of use to us. (1 Mark)</p>	<p>qualify, the two answers below are examples]</p> <ul style="list-style-type: none"> a. Walkthroughs inside a restaurant - so we can get a feel for the place b. First-person POVs of amusement rides, cruise ships, opera house seats, etc - again so we can get a preview <p>[please don't use the answers below :)] Featured snippets: Google's attempt to answer the query right on the search results page. It provides a brief summary of an answer to the search query displayed at the top of the search results. Eg. Displayed as paragraph, or list or table. (+1) Rich snippets : gives users a convenient summary information about their search results at a glance. It allows users to assess the relevance of a search result before going to the actual website. Eg. address, hours, directions, events, etc. (+1)</p>
<p>4. As you know, genAI(generative AI) is so-called because it can generate content (text, images, video, audio, more). How will this adversely affect search in the (near, even) future? Explain carefully (don't write a vague answer!). (1 Mark)</p>	<p>Search will get 'polluted' by genAI content, making the results unreliable and incorrect - eg. image-searching for 'Picasso' will pull up genAI images in the style of Picasso, in addition to genuine Picasso paintings, users will/might not know the difference.</p> <ul style="list-style-type: none"> • The search engines may not be able to distinguish between genuine and fake content. • Fake images/videos can cause spread of misinformation/disinformation if they get indexed by the search engines. • AI generated content can sound or look unnatural which can make the users feel uncomfortable or distrustful about the source. <p>(+0.5 for each effect with explanation)</p>

<p>5. ChatGPT (for example) is said to “hallucinate” sometimes (or a lot of times, depending on the type of questions) - an unfortunate term (because only minds can hallucinate!) used by companies who serve this kind of AI products. This means that the bot provides an incorrect answer(which we can verify using our own knowledge or experience or by doing a good old search!) WHAT mechanism (in the algorithm) causes this to happen? Please be specific. (1 Mark)</p>	<p>The algorithm has no understanding of what underlies language. It merely computes words one by one based on what is in its embedding, without being to assess if its output is factual or not.</p> <p>Bias : ChatGPT can provide incorrect answers to few questions because of an inherent bias caused by its limitation to understand real-world or/and lack of proper training data (imbalanced data).</p> <p>(+1 for similar explanation)</p>
<p>6. We typically write code (eg for your homeworks #2 to #5) to make use of IR algorithms. An alternative way is to use ‘nodes’ (a node is a box-like representation that encapsulates a specific task by executing task’s code) and WIRE them visually, like so –, WHAT would be two specific (and different from each other) advantages of switching to this way of working (using nodes, as opposed to coding) (1 Mark)</p>	<ul style="list-style-type: none"> • Simplicity : users without coding knowledge can also follow all the steps. • Language independent : coders from different backgrounds can interpret the flow of work and replicate it in the programming languages that they are comfortable in. • Less erroneous : syntactical errors caused by coding languages can be avoided. • Visual representation : small space can capture the flow of a program better compared to the more complex “class” structures enforced by programming. • The visual representation is ‘self documenting’, making it easy to modify by others <p>(+0.5 for each unique advantage - can be beyond the list mentioned above)</p>
<p>7. Consider the diagram below (the bottom part is simply a zoomed-in portion of the top) What algorithms did we study that results in such a collection of polygons? Why is each polygon convex? (1 Mark)</p>	<p>The figure shows a Voronoi diagram that is obtained from a k-NN classification algorithm (+0.5)</p> <p>Each item is placed in a group such that it is closer to that group than any other group, hence each polygon is convex. More: If a polygon is concave, that means that that particular item could be closer to the center of a neighboring polygon compared to the one it is a part of currently, hence nearest neighbor based polygons cannot be concave. Another explanation: the polygon’s edges are formed by pairwise perpendicular bisections</p>

	<p>of adjacent centroids, this will necessarily lead to a convex polygon. (+0.5 for similar explanation)</p>
<p>8. We revised the diagram on and off many times: Now that the course is over, how does it summarize the course? Pick 4 specific and different IR tasks we studied during the course (including “Assorted topics” pair of lectures) and explain (in a line or two) each, in terms of the three pieces of our diagram. (1 Mark)</p>	<p>Recommendation systems</p> <ul style="list-style-type: none"> - The user’s previous interests are used as an input along with the current query by the agent to give personalized content suggestions. <p>Question Answering</p> <ul style="list-style-type: none"> - The input given by the user is usually in the form of a natural language question (Who/what) and maybe specific, the agent has to identify entities and find relevant answers associating the nature of the question and the entity to provide relevant content back to the user. <p>Related search using Clustering</p> <ul style="list-style-type: none"> - The user input is the main query, the agent has to identify queries similar to the input query that may help the user get more information. This additional content may be displayed separately as “related search”. <p>Location-Based Search</p> <ul style="list-style-type: none"> - In such specific searches, one major part of the query from the user would be the location and the proximity. The agent has to identify results specific to that proximity mentioned by the user to display content. <p>LLM based agents</p> <ul style="list-style-type: none"> - The user input is usually a query in natural language. The agent must use advanced techniques to identify relevant information and *answer* using “knowledge” and present result content as a natural language response. <p>Basic (keyword) search</p> <ul style="list-style-type: none"> - the agent is a classic search engine that uses its inverted index along with ranking algorithms (including PageRank) to return a ranked sequence of results to the user.

	<p>(+0.25 for each task with relevant explanation - can be beyond this list)</p>
<p>9. How do recommendation engines work? (2 Marks)</p> <p>What are two different uses for them when we search? (2 Marks)</p>	<p>A recommendation system is any system which provides a recommendation/prediction/opinion to a user on items.</p> <ul style="list-style-type: none"> • Content-based filtering : uses item similarity/clustering to recommend items like ones you like. • Collaborative filtering : uses user similarity i.e links between users and the item they chose as the basis of recommendation. <p>(+2 for a similar explanation for content and collaborative explanation)</p> <p>Two different uses -</p> <ul style="list-style-type: none"> • Recommending content that is similar to the user's preference / promoting content may or may not be popular but similar to user taste / "sell items from long tail". • Increase diversity by recommending content that the user hasn't discovered yet. • (more generally) Recommending items bought together, nearby tourist attractions (that the user did not search for), similar songs... <p>(+1 for a each use - can be beyond this list)</p>
<p>10. What is vector similarity search?</p> <p>For what two different IR tasks are vector DBs useful? Name and explain why we cannot do it without vector DBs (1 Mark)</p>	<p>The main idea is that if data and documents are alike, their vectors will be similar. In vector space, each object or data point is represented as a vector in a high dimensional space. The goal is to find objects that are similar to a given query vector based on some similarity measure.</p> <p>(+0.5 for a similar explanation)</p> <p>IR tasks where vector DBs are useful</p> <ul style="list-style-type: none"> • Image and Video Retrieval: By representing images or video frames as vectors, vector databases enable efficient search for visually similar content

	<ul style="list-style-type: none"> • Natural Language Processing (NLP): By representing text documents or sentences as vectors, vector databases allow efficient similarity search based on semantic or contextual similarities • Recommendation Systems: By representing items as vectors based on their attributes or features, vector databases enable efficient retrieval of similar items <p>(0.5 for any two of the above two answers with correct explanation)</p>
<p>11. Summarize any two HWs #2 through #5 - what algorithm was used and what task did it help accomplish?</p> <p>(1 Mark)</p>	<p>HW2: Crawler4j, to develop a develop a 'crawler' to 'crawl' (visit, fetch contents of) a list of URLs,</p> <p>HW3: Uses Map Reduce to build an inverted index of a data</p> <p>HW4: Uses Lunr and Solr for indexing and searching documents</p> <p>HW5: Uses Weaviate to set up, load and query vectorized data</p> <p>(+1 for any two with correct explanation)</p>
<p>12. What does OPL stand for, in the OPL stack?</p> <p>What is its main use?</p> <p>(1 mark)</p>	<p>"OPL" is the name we could give to such a stack comprised of O(penAI)+P(inecone)+L(angChain)</p> <p>(+0.5 for the correct full form)</p> <p>It has increasingly become the industry solution to overcome the two limitations of LLMs:</p> <ol style="list-style-type: none"> 1. LLMs hallucination: chatGPT will sometimes provide wrong answers with overconfidence. One of the underlying causes is that those language models are trained to predict the next word very effectively, or the next token to be precise. Given an input text, chatGPT will return words with high probability, which

	<p>doesn't mean that chatGPT has reasoning ability.</p> <p>2. Less up-to-date knowledge: chatGPT's training data is limited to internet data prior to Sep 2021. Therefore, it will produce less desirable answers if your questions are about recent trends or topics.</p> <p>(+0.5 for the any two valid reasons with correct explanation)</p>
<p>13. How do RDF triples make search better? (1 mark)</p>	<ul style="list-style-type: none"> - RDF triples enable semantic search, which goes beyond simple keyword matching - RDF triples facilitate the integration of data from multiple sources through the use of unique identifiers (URIs) and linking - DF triples enable the expansion and enrichment of search queries - RDF triples provide a structured and standardized way to represent data - RDF triples can be used to represent ontologies and knowledge graphs, which enable reasoning and inference <p>(+1 for any of the above two reasons with correct explanation)</p>
<p>14. Name four algorithms we looked at, for IR tasks, that reply to iteration or recursion. for each briefly explain how iteration or recursion helps? (4 Marks)</p>	<ul style="list-style-type: none"> - Page Rank Algorithm: iteration is used to improve the accuracy and convergence of the ranking calculation - snippet generation algorithm: iteration plays a crucial role in generating concise and relevant snippets from a larger text document - Basic spelling correction algorithm uses iteration/recursion to improve the accuracy of the correction process. - Google's Query Processing Algorithm - K-Means - Hierarchical clustering algorithm - Rocchio Algorithm - KNN <p>(+1 for a correct a algorithm with correct explanation)</p>

	<p>(Note: The answer is not limited to the above points, any other algorithms taught in class are valid)</p>
<p>15. Name and briefly discuss 4 ML based algorithms we looked in class</p> <p>(4 Marks)</p>	<ul style="list-style-type: none"> - K-Means - Hierarchical clustering - Rocchio Algorithm [not ML :)] - KNN - Neural networks (eg in Transformers that are in (chat)GPT) <p>(+1 for a correct a algorithm with correct explanation)</p> <p>(Note: The answer is not limited to the above points, any other algorithms taught in class are valid)</p>
<p>15. Consider the image, pick four algorithms that are useful in IR. Explain the algorithm and how they are useful in IR.</p> <p>(4 marks)</p>	<ol style="list-style-type: none"> 1. Naive Bayes: <ul style="list-style-type: none"> • Text categorization • Sentiment analysis 2. Support Vector Machines (SVM): <ul style="list-style-type: none"> • Document classification • Information extraction 3. Neural Networks: <ul style="list-style-type: none"> • Convolutional neural networks (CNNs): <ul style="list-style-type: none"> • Document classification • Image retrieval • Recurrent neural networks (RNNs): <ul style="list-style-type: none"> • Text generation • Machine translation • Question answering 4. Decision Trees: <ul style="list-style-type: none"> • Document classification • Feature selection 5. Random Forests: <ul style="list-style-type: none"> • Document classification • Sentiment analysis 6. Gradient Boosting: <ul style="list-style-type: none"> • Ranking and relevance prediction 7. K-Means or Hierarchical clustering <ul style="list-style-type: none"> • document clustering • query expansion • text expansion 8. KNN - for classification <p>(+1 for one correct algorithm with correct explanation)</p>

	<p>(Note: The answer is not limited to the above points, any other algorithms mentioned in the image that has a use case in IR are valid)</p>
<p>17. Tiktok's recommendation engine uses a specific data structure to optimize it. What is the name of the data structure? How does it work?</p> <p>(1 mark)</p>	<p>The data structure is a 'cuckoo hash'. It works by using two hash tables/functions rather than just one so that if a collision occurs in an incoming key, the existing key is kicked out to a new location in the 'other' table. We can accept 'loose' explanations as long as the use of two tables (or functions) is mentioned.</p> <p>TikTok uses Monolith DS to optimize its recommendation system.</p> <p>(+0.5 for the correct DS)</p> <p>(+0.5 for correct explanation)</p>
<p>18. given two vectors, what two similarity measures can we calculate?</p> <p>What do vectors have to do with IR?</p> <p>(1 mark)</p>	<p>Similarities we can calculate:</p> <ul style="list-style-type: none"> - Cosine similarity - Euclidean Distance <p>(0.5 for any two similarity measures taught in class)</p> <p>IR often involves representing documents, queries, or data points as vectors in a high-dimensional space. By representing textual or numerical features of documents as vectors, various similarity measures can be applied to compare and rank documents based on their similarity to a given query</p> <p>(+0.5 for a similar explanation)</p>