

Data Science Canvas		Project:	Research Paper Recommendation Service				
		Team:	ScholarStream				
Problem Statement			Execution & Evaluation		Data Collection & Preparation		
<b>Business Case &amp; Value Added</b> There is explosion of the research articles, and it is difficult to keep up with the newly introduced concepts and trending topics.	<b>Model Selection</b> Embedding models are known to model a vector space that encodes real world concepts and therefore help with semantic similarity.  We will use, TF-IDF as baseline 1, all-minilm-l6-v2 as baseline 2	<b>Model Requirements</b> <ul style="list-style-type: none"> <li>Semantic understanding.</li> <li>Domain relevance</li> <li>Scalability</li> <li>Low latency</li> <li>Input flexibility</li> <li>Low cost</li> </ul>	<b>Skills</b> What skills are needed to provide the data and model development? <ul style="list-style-type: none"> <li>EDA</li> <li>Developing application using python framework like FASTAPI</li> <li>Vector DB</li> <li>NLP</li> <li>Machine Learning</li> <li>Containerization</li> <li>Cloud Deployment</li> </ul>	<b>Model Evaluation</b> <ul style="list-style-type: none"> <li>Precision@K</li> <li>Recall@K</li> <li>MRR</li> <li>nDCG@K</li> </ul> Real time monitoring: <ul style="list-style-type: none"> <li>Latency</li> <li>Throughput</li> <li>Error rate</li> <li>Resource usage</li> </ul>	<b>Data Storytelling</b> A simple ui with detailed research paper information, recommendations and similarity scores.	<b>Data Selection &amp; Cleansing</b> The data is readily available in arxiv. We filter out papers that have non ascii characters in the title or abstract.	<b>Data Collection</b> No additional data is required.
<b>Data Landscape</b> Yes data is already available from arxiv Additional data can be generated by extracting pdf content for certain categories. But, for the purpose of this project, we will limit ourself to just the metadata (title, abstract, authors etc.)		<b>Software &amp; Libraries</b> <ul style="list-style-type: none"> <li>Numpy</li> <li>Pandas</li> <li>Sci-kit learn</li> <li>Vector db like milvus db</li> <li>LLMs (ex gemini-flash-2.5)</li> <li>Langchain</li> <li>HuggingFace</li> <li>FastAPI</li> </ul>			<b>Data Integration</b> Data integration is not required as it is from a single source	<b>Explorative Data Analysis</b> We only use papers from 'cs.ai' category. The data is clean and doesn't require EDA.	