

Material data extraction from journal articles using large language models



Aravindan Kamatchi Sundaram, Devathi Sai Mani Kumar, Rohit Batra

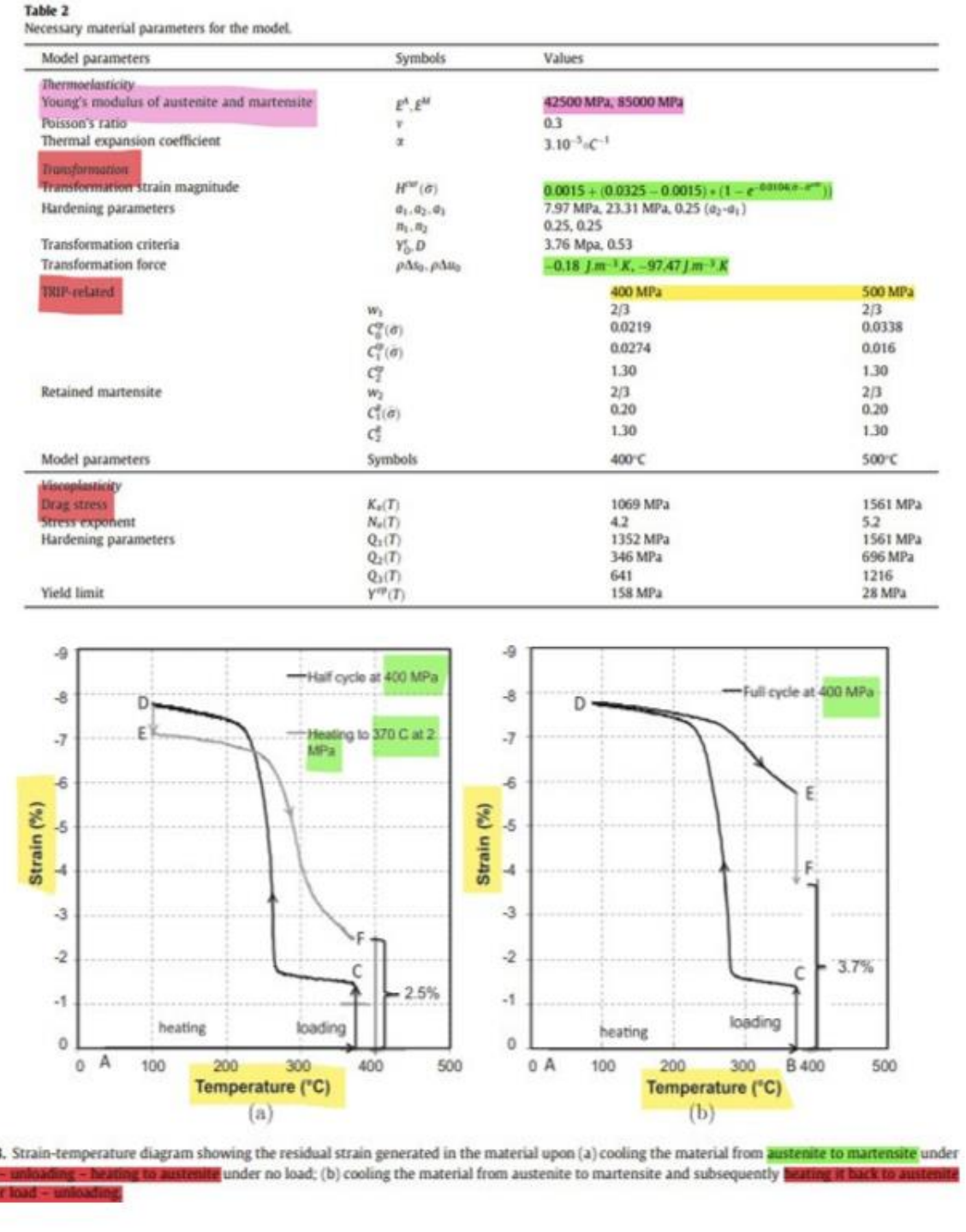
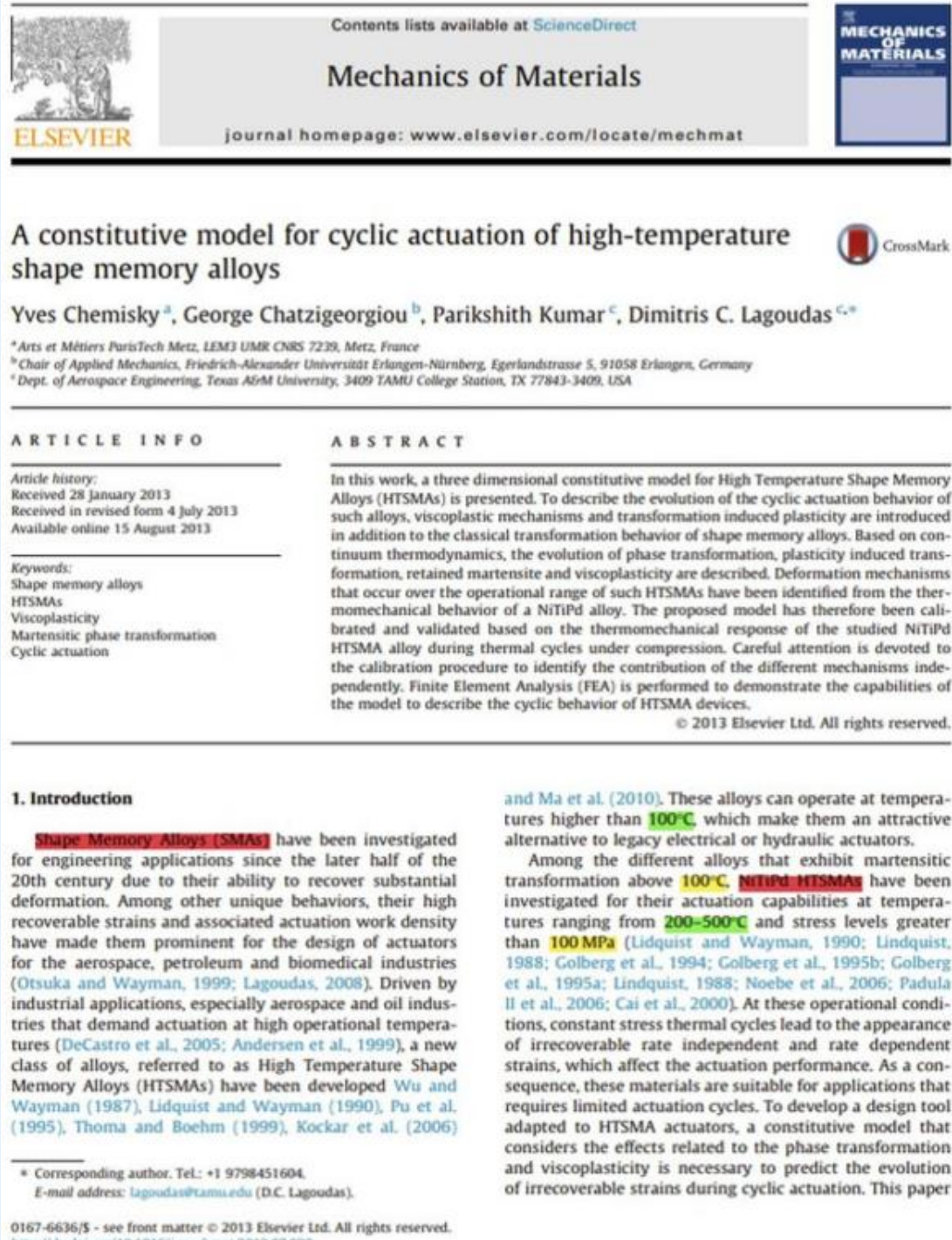
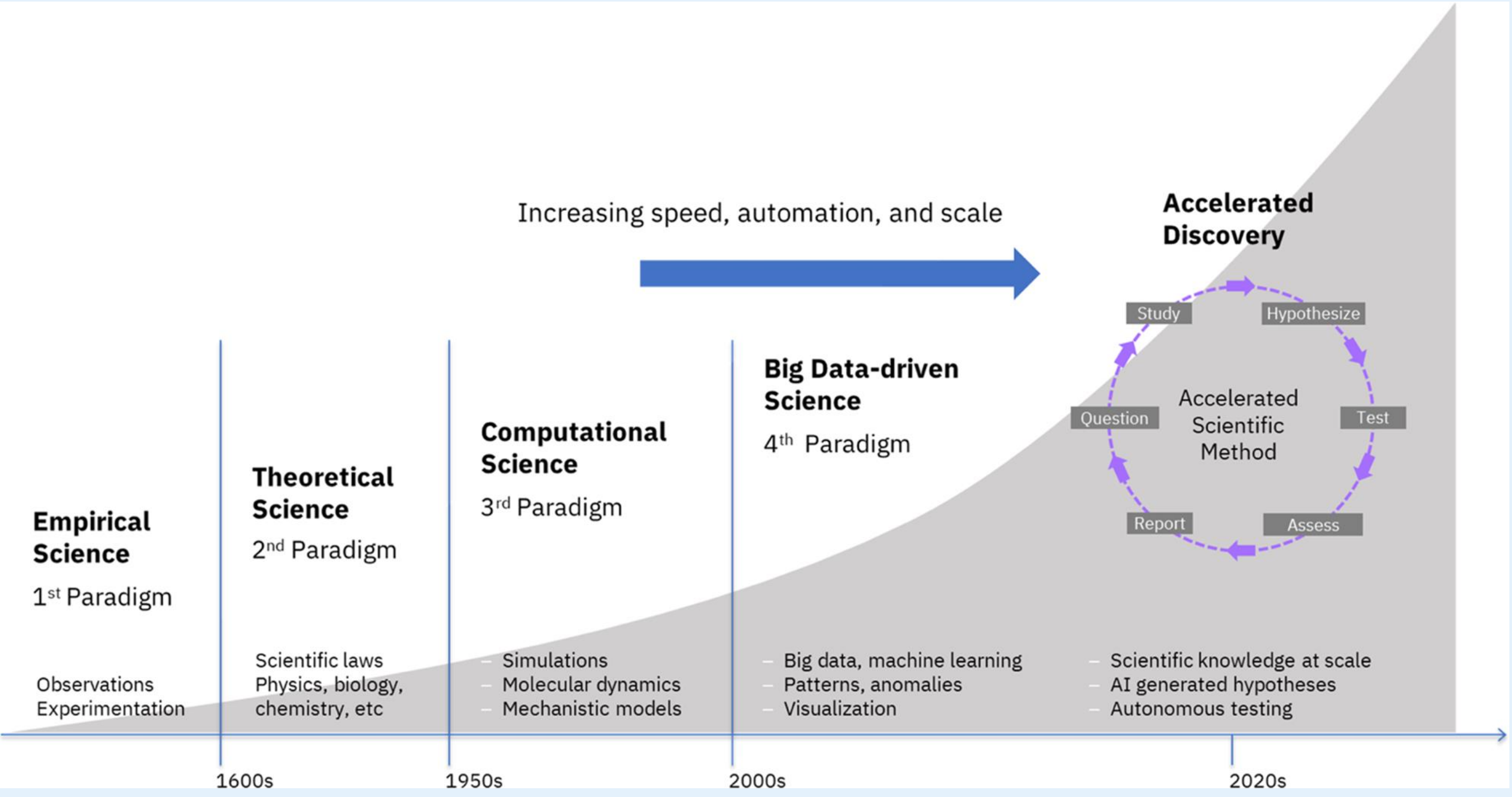
Department of Metallurgical and Materials Engineering

Indian Institute of Technology Madras, Chennai

Objectives

- To establish a framework to extract materials' data in a structured format from research sinks.
- To test the framework on a sub-field in materials science and evaluate the performance of the framework when used with different LLMs.

Motivation



Materials Data in academic works

Methods

0. Data Collection:

Elsevier API+ XML Parsing of research articles to obtain the different sections of the paper.

1. Prompt Engineering:

Prompts: Context + few shot examples+ formatting instructions+ article

Context:

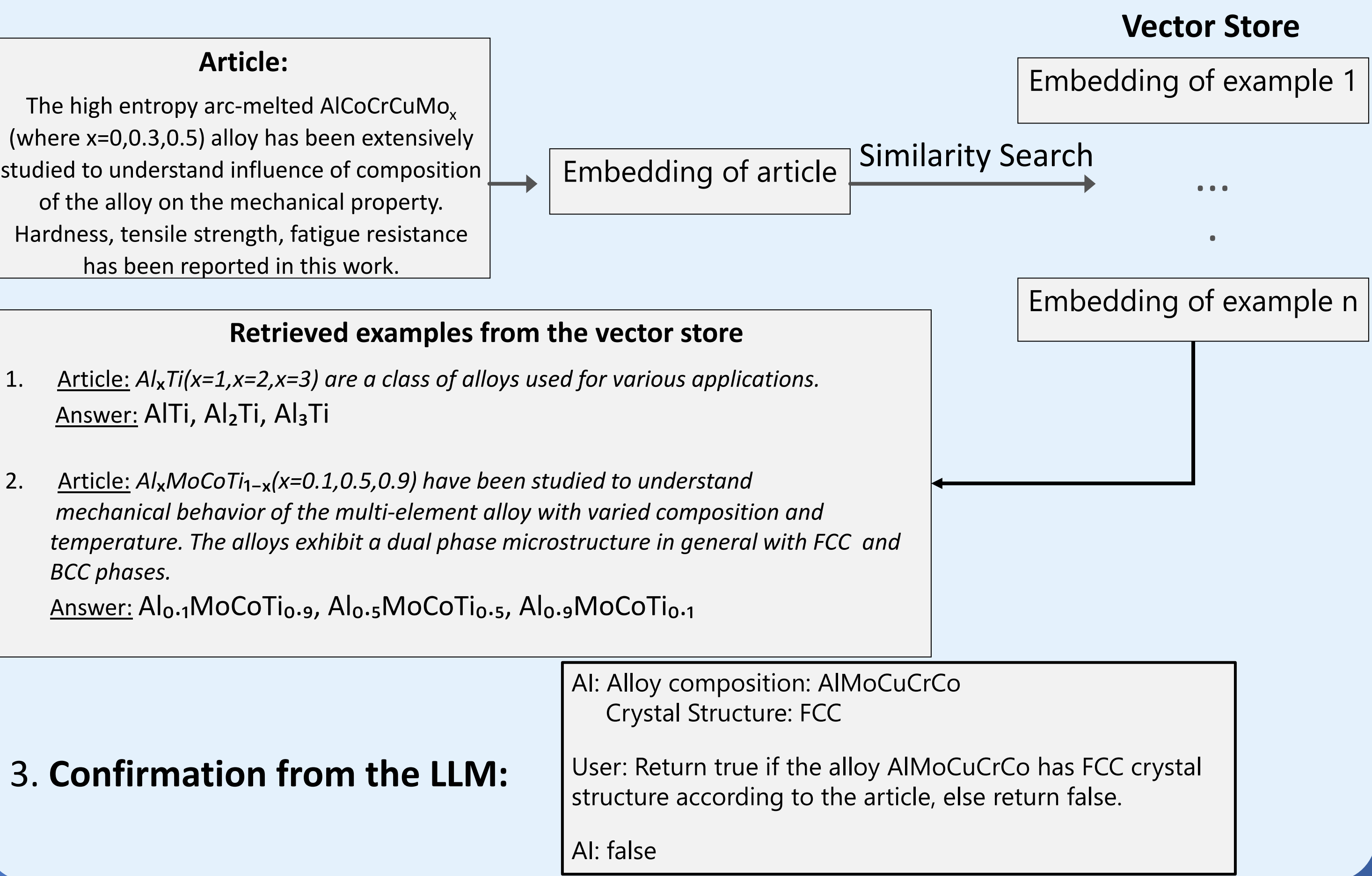
- An alloy is a solid mixture of multiple elements. Its chemical composition conveys the amount of each element present in the alloy. For example, an alloy with chemical composition AlCoCr contains elements Al, Co and Cr in equal amount. Similarly, an alloy with chemical composition $Al_2Co_3Ni_5$ contains 2 parts Al, 3 parts Co and 5 parts Ni.
- Here are some of the common structures adopted by alloys
- ...
- Abbreviations for subsequent thermo-mechanical processes are: CR=Cold Rolled; FC=Furnace Cooled; FOR=Forged; HIP=Hot Isostatic Pressing; HPT=High Pressure Torsion; HR=Hot Rolled; VHP=Vacuum Hot Pressed; WQ=Water Quenched

Few shot examples:

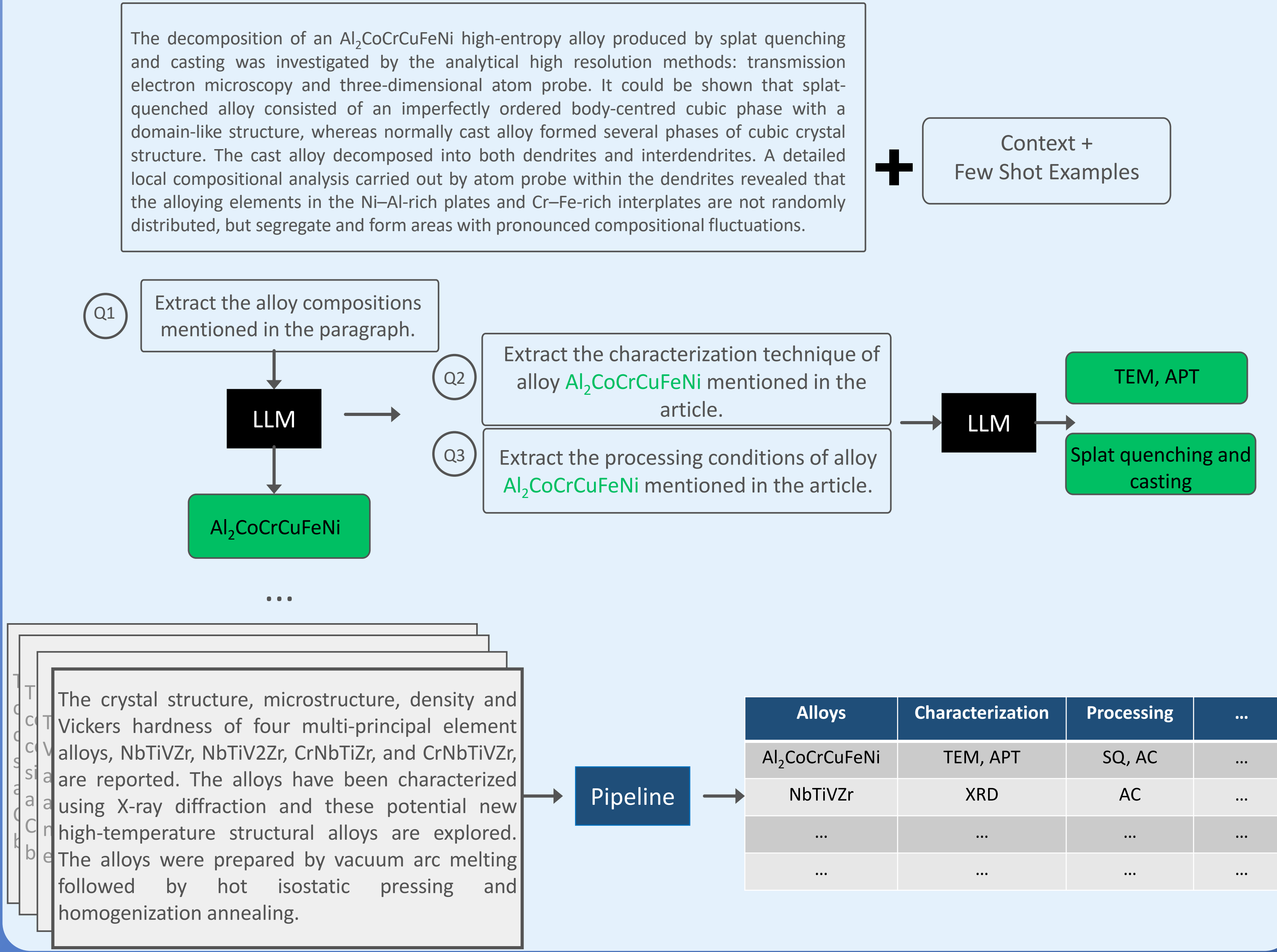
- Article: We have experimented on the chemical and mechanical high-entropy alloys CoCrFeMnNi. We have also done experiments with alloys after substituting elements, Ti for Co Mo or V for Cr. The factors affecting stability of various phases is studied. Answer: CoCrFeMnNi, TiCrFeMnNi, CoMoFeMnNi, CoVFeMnNi
- Article: AlCrFeCoNiCuTi alloy contains BCC1 phase, BCC2 phase and a FCC phase. AlCrFeCoNiCuV alloy contains two different phases with BCC, FCC structure, respectively. Answer: AlCrFeCoNiCuTi, AlCrFeCoNiCuV
- ...
- Article: $Al_xTi(x=1, x=2, x=3)$ are a class of alloys used for various applications. Answer: AlTi, Al_2Ti , Al_3Ti

2. Retrieval Augmented Generation:

RAG retrieves similar examples from a pre-written collection of domain specific examples.



Data extraction workflow



Results and discussions

| | Alloys | Characterization techniques | Processing conditions | Properties | Value | Units |
|-----|---|--|---|---------------------------|-------|-------------------|
| 1 | NbMoTaW | X-Ray diffraction | Vacuum arc-melting | Yield Stress | 1390 | MPa |
| 2 | NbTiVZr | X-Ray diffraction | As-cast, Splat quenched | - | - | - |
| 3 | Al _{0.5} CoCrCuFeNi | Scanning electron microscopy, Electron dispersive spectroscopy | Arc-melting, water quenched and cold rolled | Hardness | 208 | HV |
| 4 | Ti _{0.5} CrFeCoNiAl _{0.75} Cu _{0.25} | Scanning electron microscopy | Arc-melting | Compressive strength | 2697 | MPa |
| 5 | Ti _{0.5} CrFeCoNiCu | Scanning electron microscopy | Arc-melting | Compressive strength | 3135 | MPa |
| 6 | Al _{0.3} CoCr ₂ FeNi | Scanning electron microscopy, X-Ray diffraction | Arc-melted and water quenched | Vickers Hardness | 343 | HV |
| 7 | CoCrFeMnNi | X-Ray diffraction | Arc melting | - | - | - |
| 8 | NbTiV ₂ Zr | Pyconometry, Scanning electron microscopy | Vacuum arc melting, hot isostatic pressing | Density | 6340 | kg/m ³ |
| 9 | Al _{0.7} CoCrFeNi | Scanning electron microscopy | Vacuum arc melting, hot isostatic pressing | Ultimate tensile strength | 1400 | MPa |
| ... | ... | ... | ... | ... | ... | ... |
| n | Al _{0.2} Co _{1.5} CrFeNi _{1.5} Ti | SEM, EDS, XRD | Arc-melting, aging | Wear resistance | 5500 | m/mm ³ |

Data extracted from academic works of multi-element alloys

- To evaluate our model, we compute:

$$\text{Recall} = \frac{|\{\text{relevant entries}\} \cap \{\text{retrieved entries}\}|}{|\{\text{retrieved entries}\}|}$$
$$\text{Precision} = \frac{|\{\text{relevant entries}\} \cap \{\text{retrieved entries}\}|}{|\{\text{relevant entries}\}|}$$

$$\text{F1-Score} = \text{Harmonic_mean}(\text{Precision}, \text{Recall})$$

- We tested our framework on papers on multi-element alloys

| MODELS: | GPT-3.5 turbo | GPT-3.5 turbo with confirmation | GPT-3.5 turbo with RAG and confirmation | GPT-4o with RAG and confirmation | GPT-4o with RAG and without confirmation | GPT-4o mini with RAG and without confirmation |
|-----------------------------|---------------|---------------------------------|---|----------------------------------|--|---|
| Alloys | 0.70 | 0.90 | 0.87 | 0.92 | 0.96 | 0.92 |
| Processing Conditions | 0.80 | 0.87 | 0.84 | 0.91 | 0.96 | 0.81 |
| Characterization techniques | 0.63 | 0.88 | 0.83 | 0.82 | 0.99 | 0.93 |
| Properties | 0.86 | 0.83 | 0.89 | 0.75 | 0.89 | 0.94 |

F1 Scores for our experiments with various LLMs on works from multi-element alloys

| | Alloys | Processing Conditions | Characterization techniques | Properties | | Alloys | Processing Conditions | Characterization techniques | Properties |
|-----------|--------|-----------------------|-----------------------------|------------|-----------|--------|-----------------------|-----------------------------|------------|
| Precision | 1 | 1 | 1 | 0.88 | Precision | 1 | 0.77 | 0.93 | 0.95 |
| Recall | 0.93 | 0.92 | 0.98 | 0.90 | Recall | 0.86 | 0.85 | 0.93 | 0.93 |
| F1-Score | 0.96 | 0.96 | 0.99 | 0.89 | F1-Score | 0.92 | 0.81 | 0.93 | 0.94 |

GPT-4o mini with RAG Analysis of the best models for the multi-element alloy data extraction task

| MODELS: | GPT-3.5 turbo | GPT4o-mini | GPT4o |
|----------------|---------------|------------|--------|
| Cost per paper | \$0.03 | \$0.003 | \$0.15 |

- Data extraction using Openai's LLMs take around 30 seconds per article.
- Automated data extraction takes significantly lesser time and performs reliably.

Future work

- Extraction of materials' data from a collection of 18000 academic works on multi-element alloys.
- Data extraction from figures and other forms of data present in the papers.
- Exploration of RAG for context retrieval, and exploration of fine-tuning to improve performance.

Acknowledgements

I would like to thank Dr. Rohit Batra for his support and my colleague Saimani.

