

MM 226 – Materials Informatics

Assignment 2 Report

Name: Samridhi

Roll No.: 230005042

Prompt Used for Each AI

I used the following prompt across four different AI models – **ChatGPT, Claude, Google Gemini, and Perplexity:**

"Based on the Excel file uploaded, you have to generate a table containing 60 data points for duplex steel containing the following:

- (1) Chemical composition of alloy
- (2) Processing condition (Heat treatments at specific case, Test temperature)
- (3) Microstructural parameters (grain size, phases, dislocation density)
- (4) Mechanical properties with units (YS, UTS, Hardness, Elongation, strain rate, strain hardening exponent and coefficient)

Deliberately make variability, missing values, and minor inconsistencies.

Convert that data into an .xlsx file."

Each AI produced an individual **.xlsx** file except Perplexity, which returned code to create the file. These were combined into a single **.xlsx** file after standardizing column names and formatting. Each sheet was converted to **.csv** for use in Python (Google Colab) for data cleaning and analysis.

Data Cleaning Methodology and steps

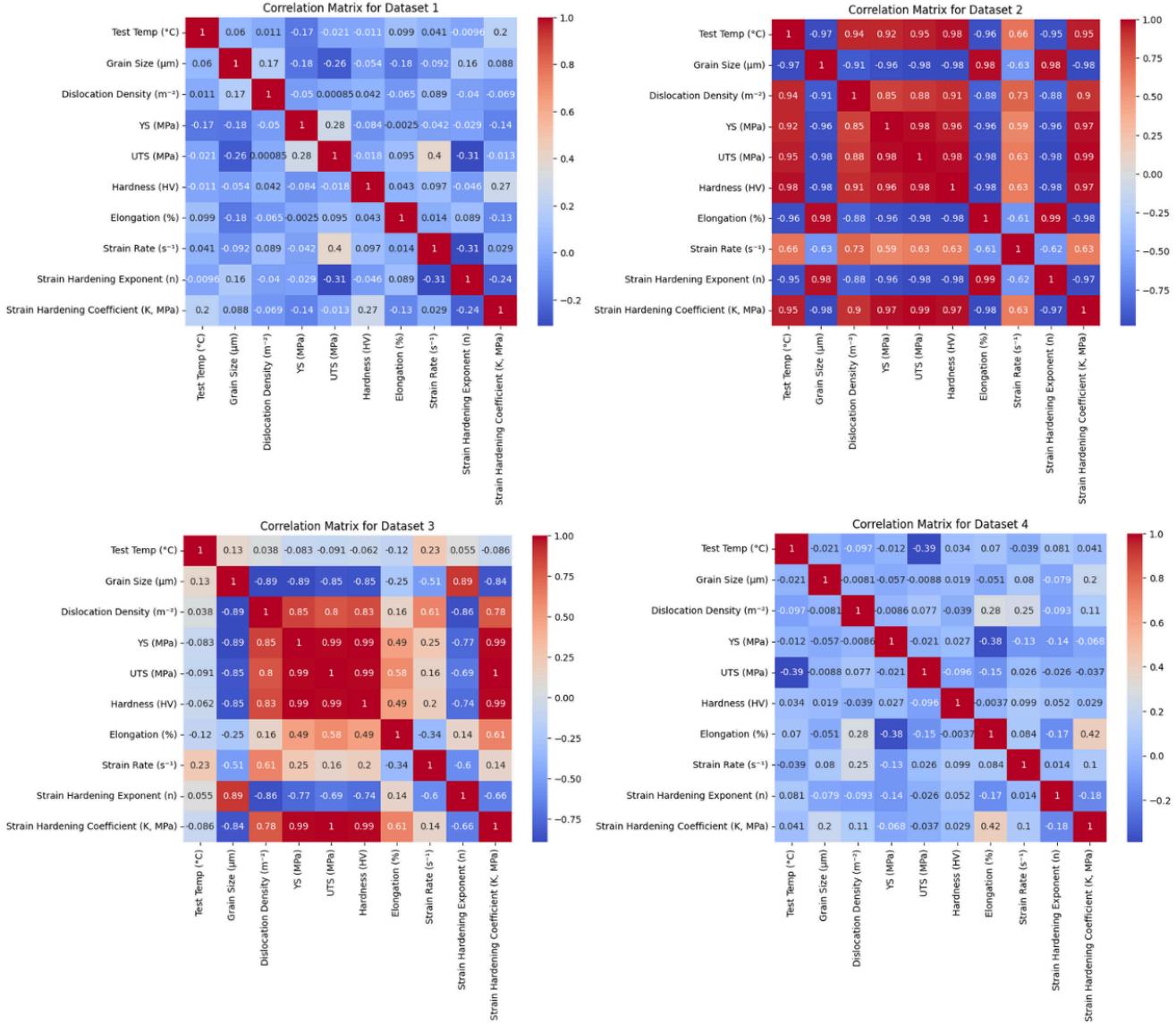
I wrote a Python code on Google Colab to clean and preprocess the dataset. The key steps involved are:

Handling Missing Values

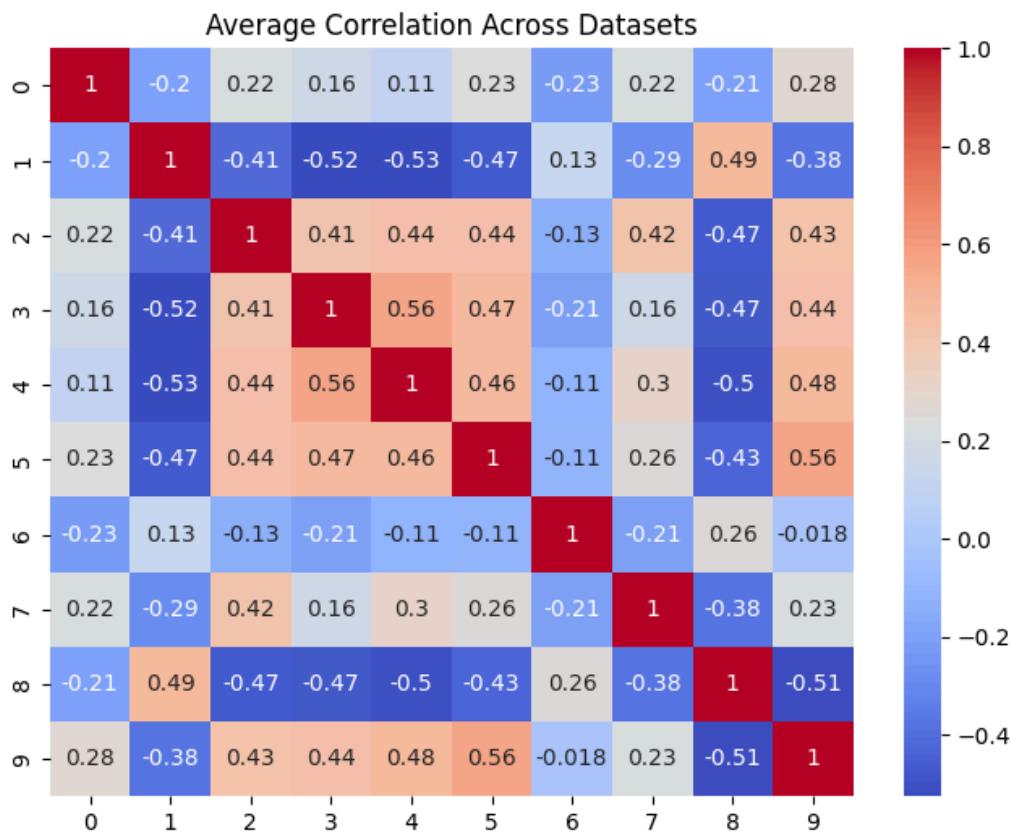
- **Test Temperature:** Replaced missing values with 25°C as this is a common standard testing temperature.
- **Phases:** Missing phase data was filled with “Austenite + Ferrite”, the typical structure in duplex steels.
- **Strain Rate:** Filled with 0.001 s⁻¹, a common quasi-static testing strain rate.
- **Other Parameters (Grain Size, Dislocation Density, YS, UTS, Hardness, Elongation, n, K):** Missing values were replaced with the mean of the respective columns.

Comparative Analysis of Four Datasets

The four datasets were compared using a range of visualizations including heat maps and pair plots and a combined heat map was also generated for high-level comparison.



Aspect	ChatGPT	Claude	Gemini	Perplexity
Data plausibility	High	High	Medium	Medium
Value ranges	Broad and realistic	Broad	Realistic but narrower	Slightly inconsistent
Inconsistencies	Minor, plausible	Minor	Some structural errors	Few invalid entries
Property relationships	Clear	Mostly clear	Noisy	Slight deviation

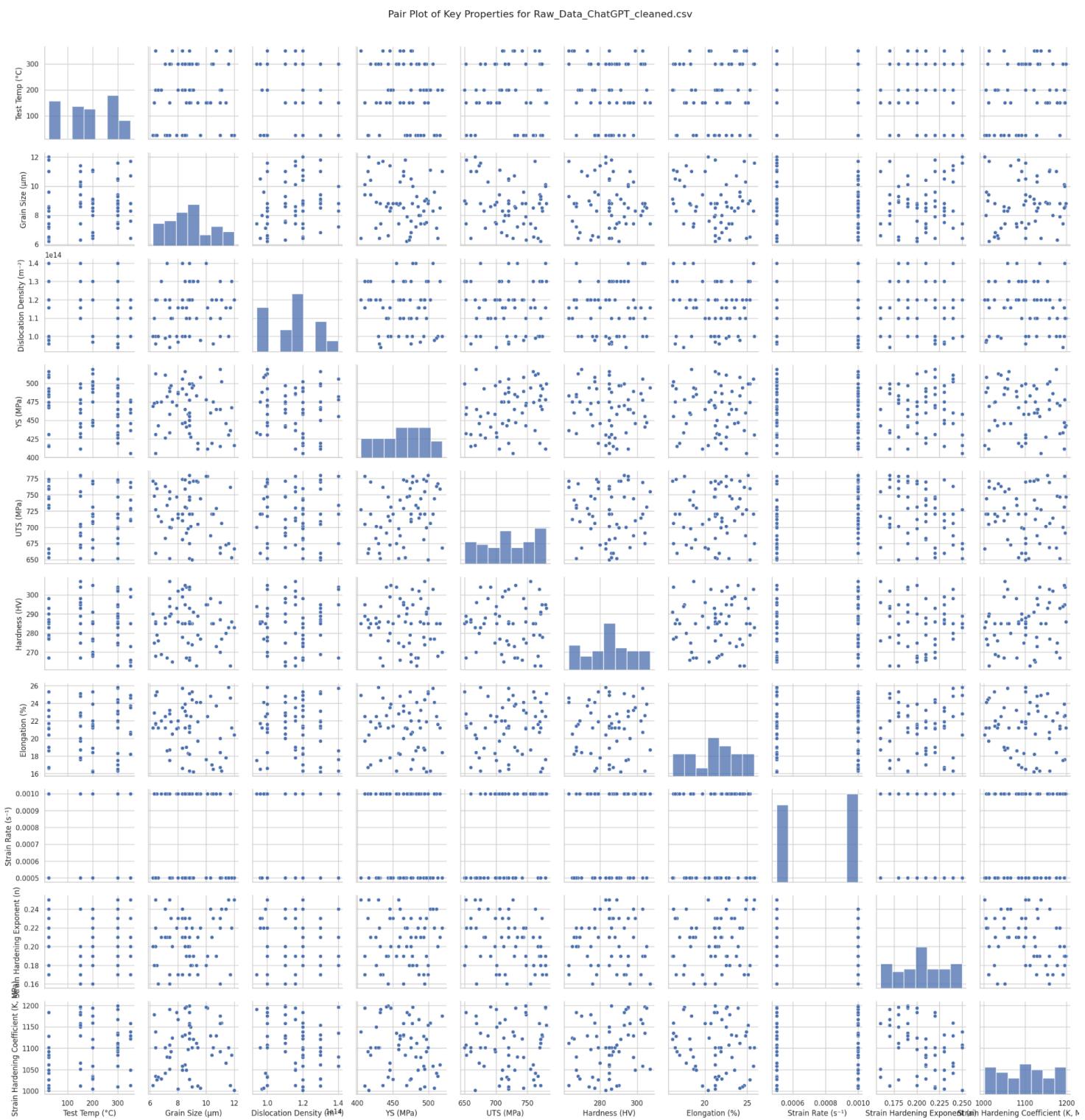


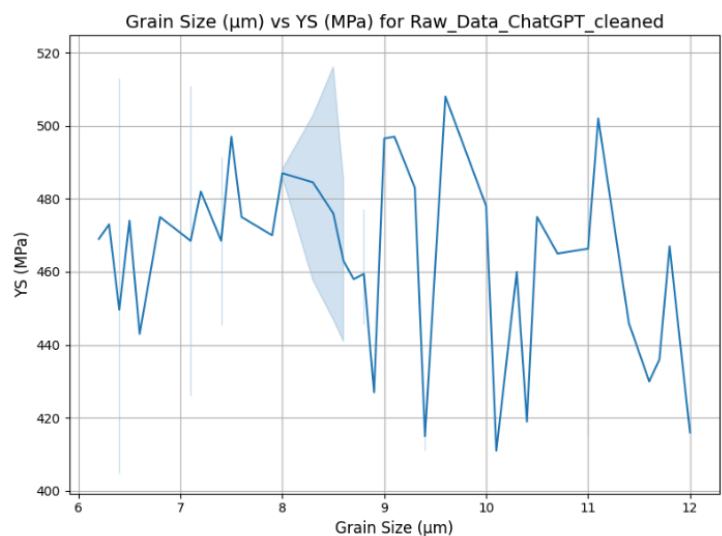
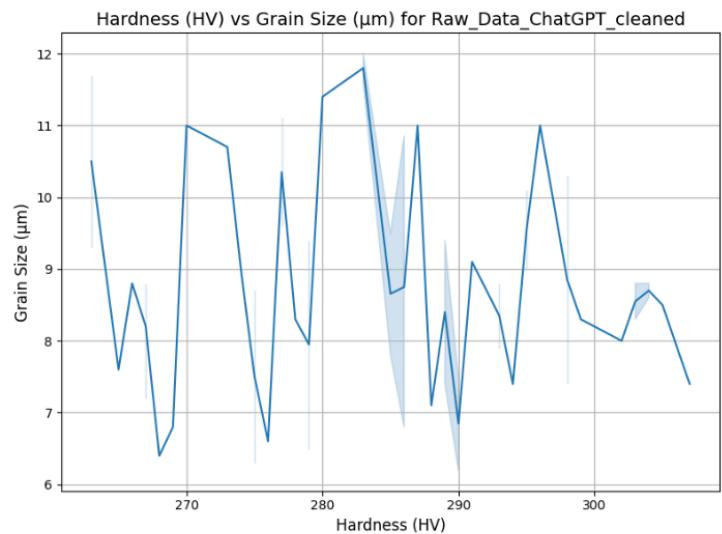
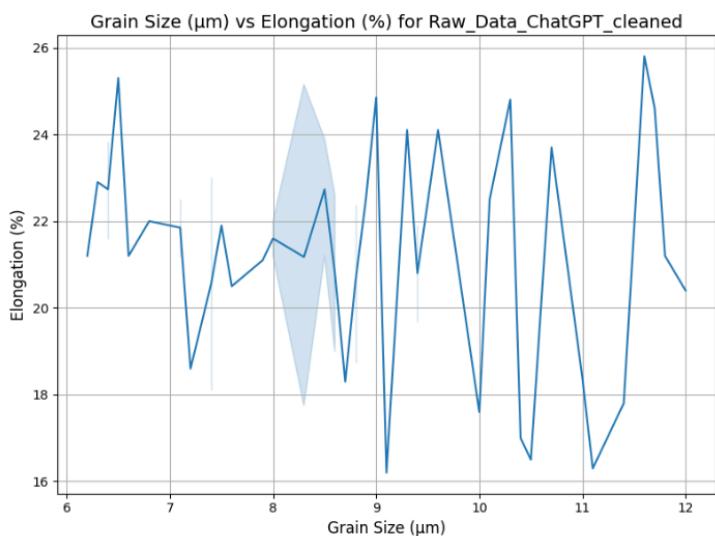
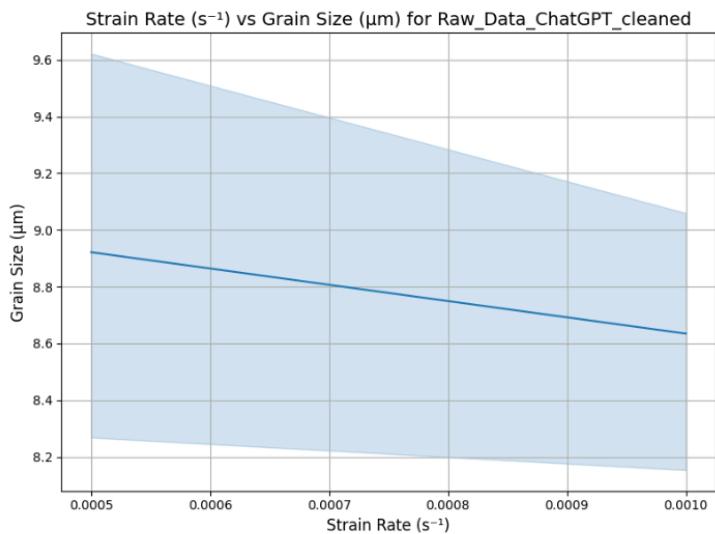
Detailed Observations:

- ChatGPT and Claude produced the most plausible and consistent datasets.
- Gemini's values were within reasonable range but had more formatting issues.
- Perplexity gave code-based output; though slightly noisier, data was useful after cleaning.

Pair Plot and x vs y Plots for Each Al

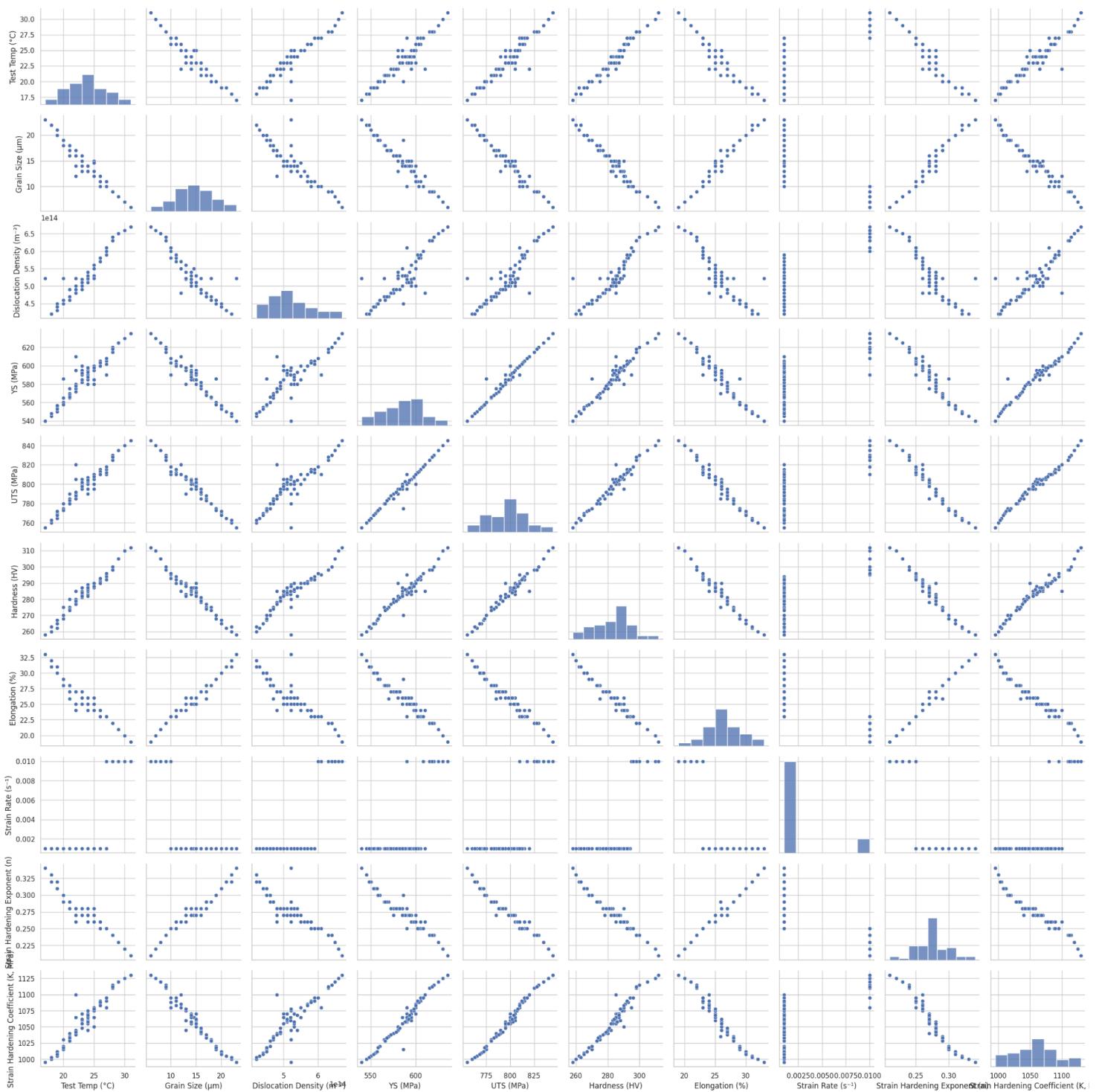
ChatGPT: The data was well-distributed with consistent trends. Strong correlation observed between YS and UTS. Slightly tighter ranges than others.

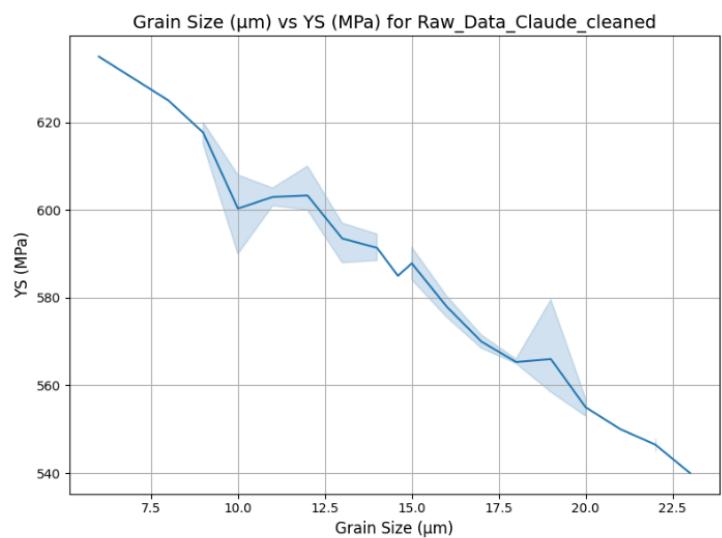
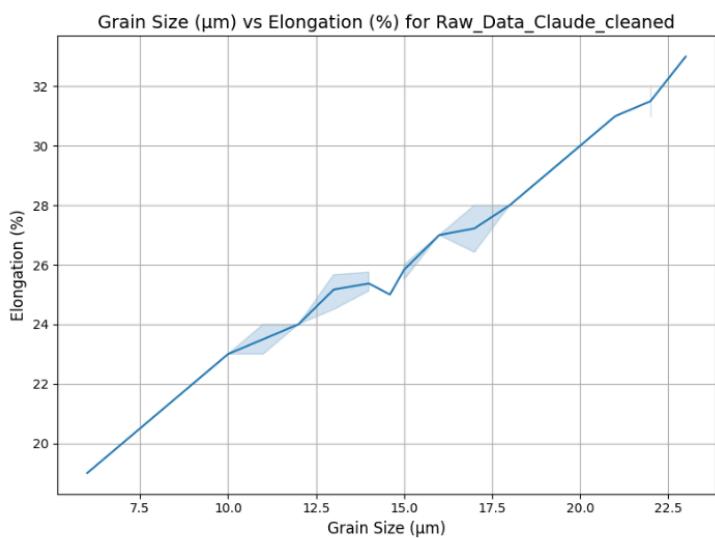
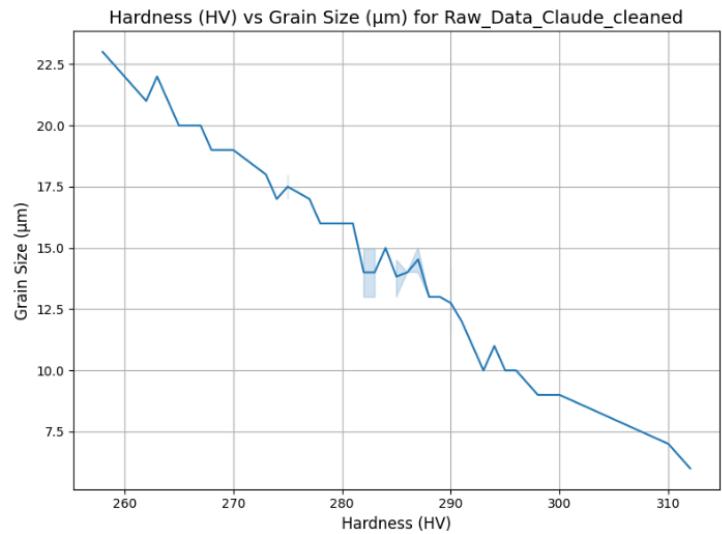
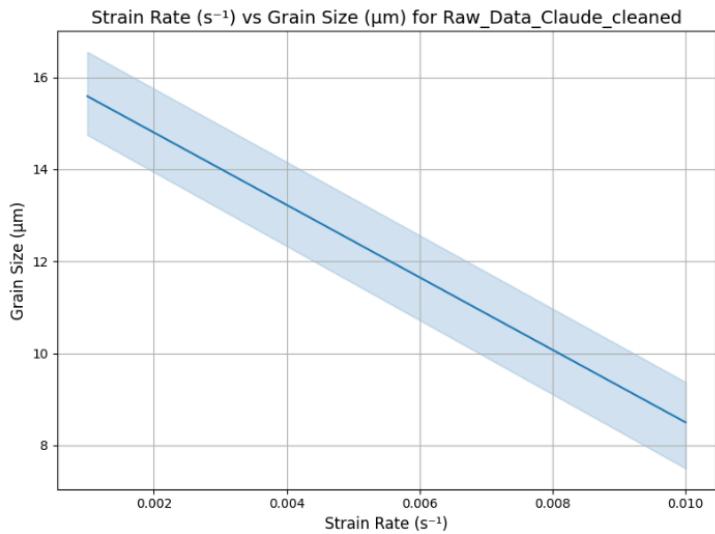




Claude: Shows good spread in grain size and hardness. Moderate noise in strain rate and elongation.

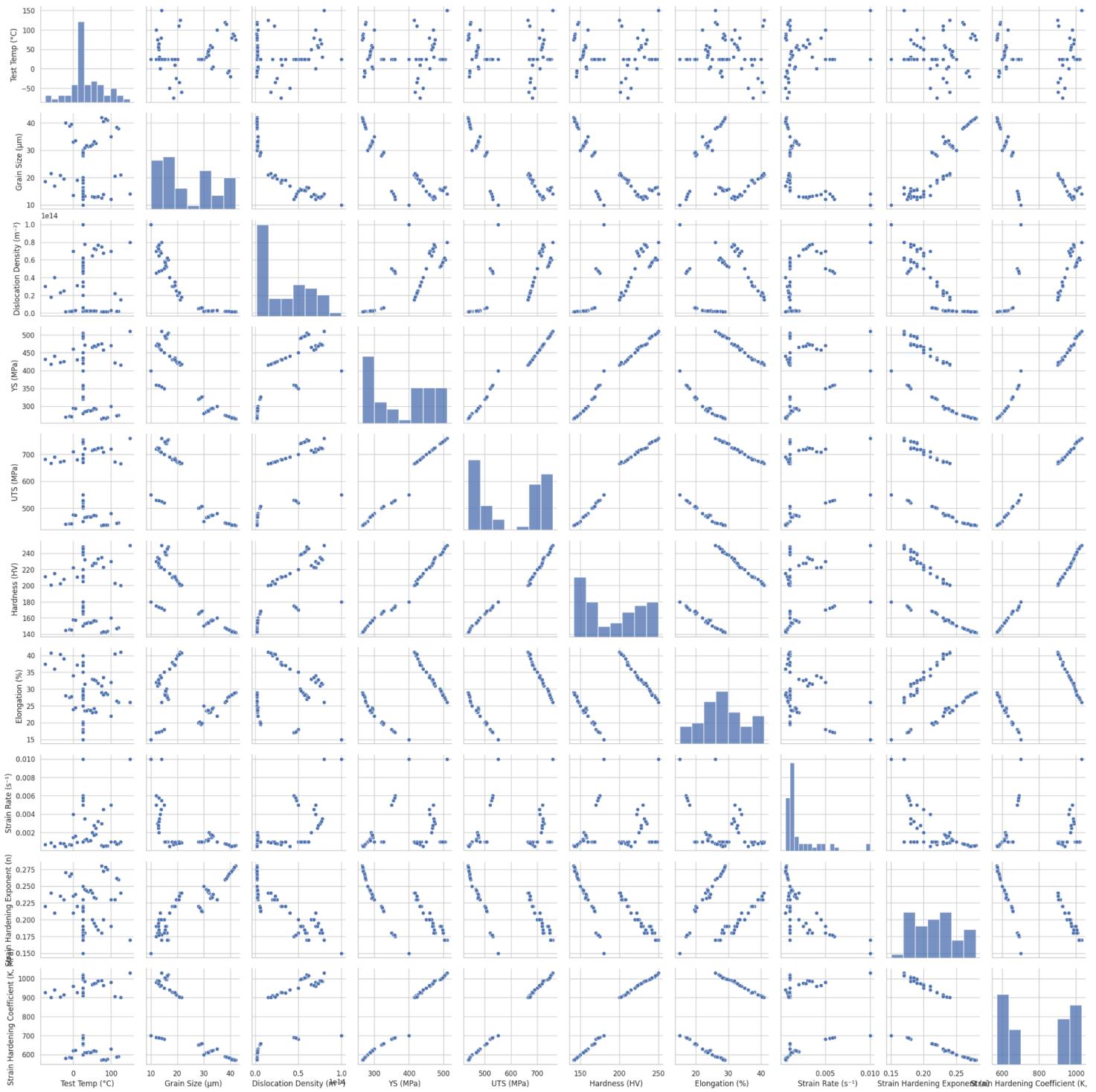
Pair Plot of Key Properties for Raw_Data_Claude_cleaned.csv

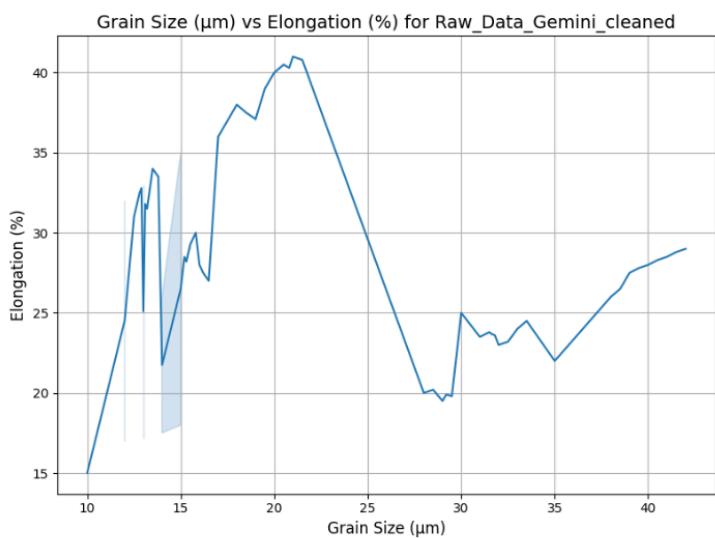
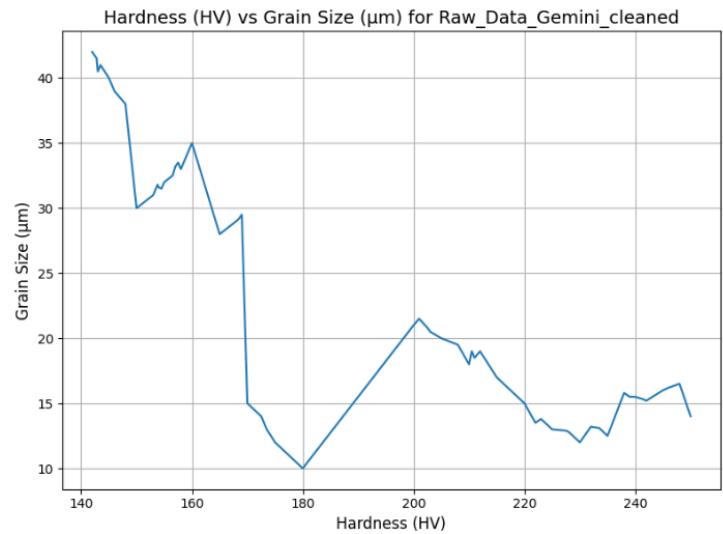
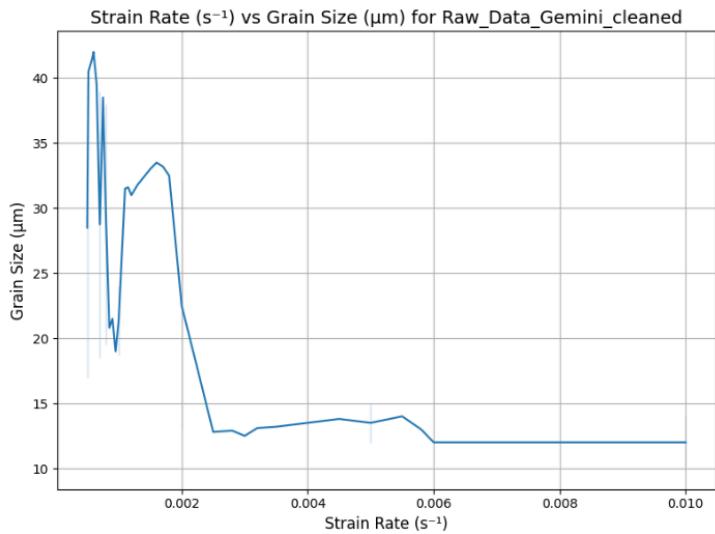




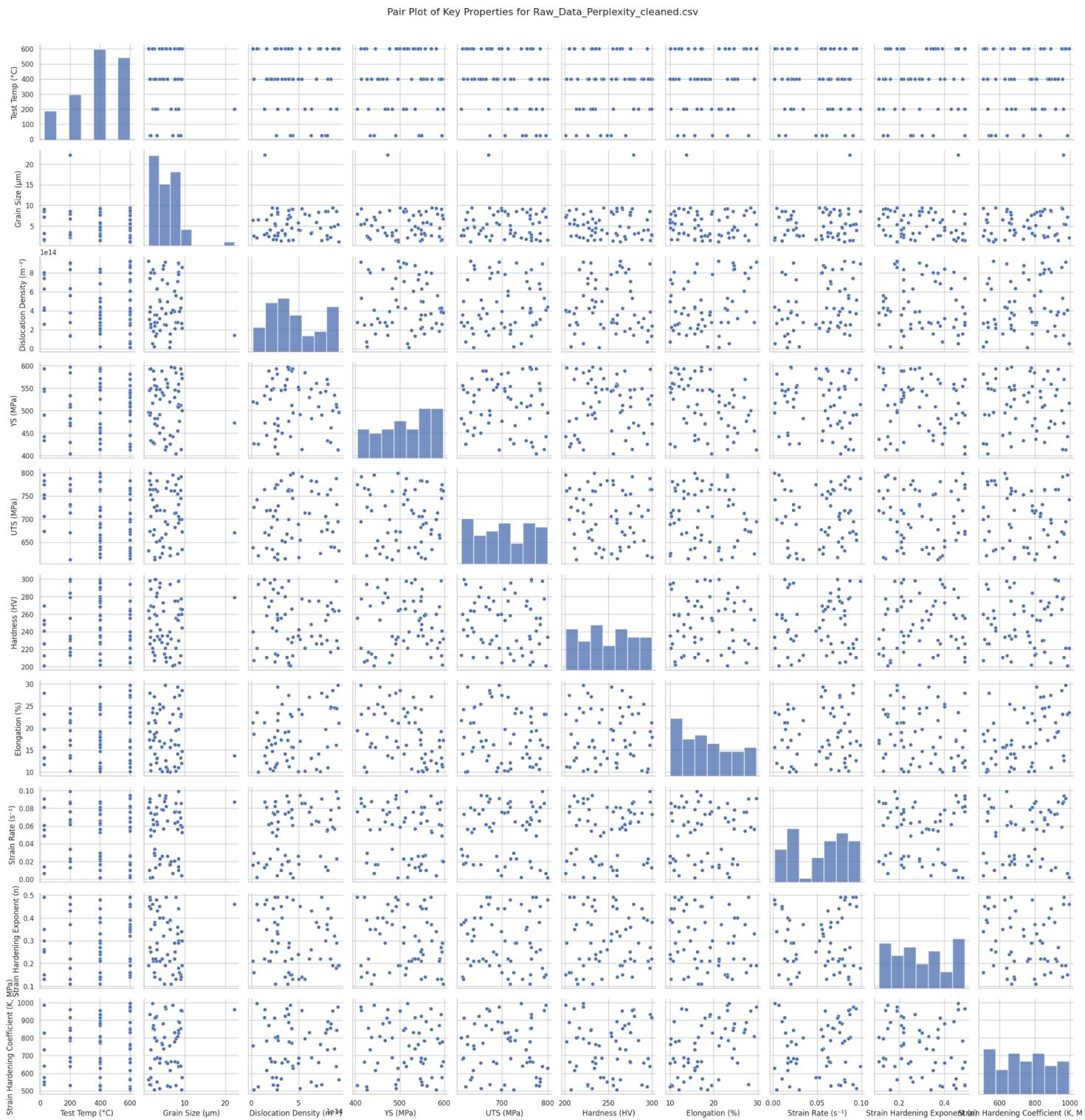
Gemini: Broader spread in microstructural features. Some missing values affected correlations. A few outlier points present.

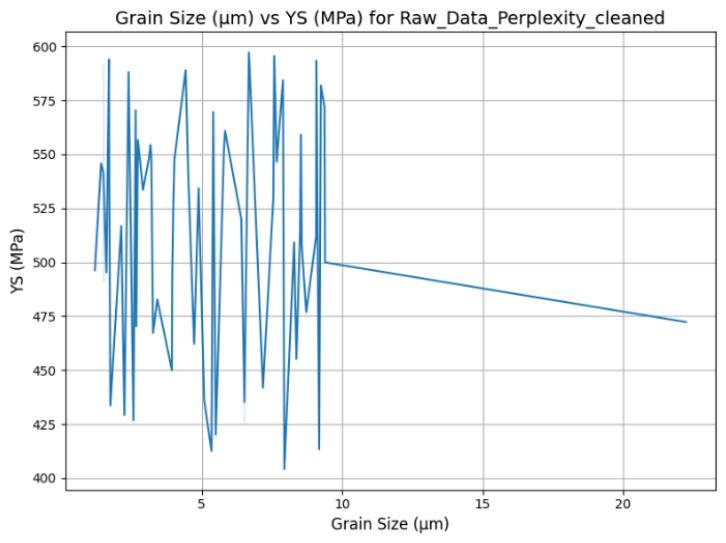
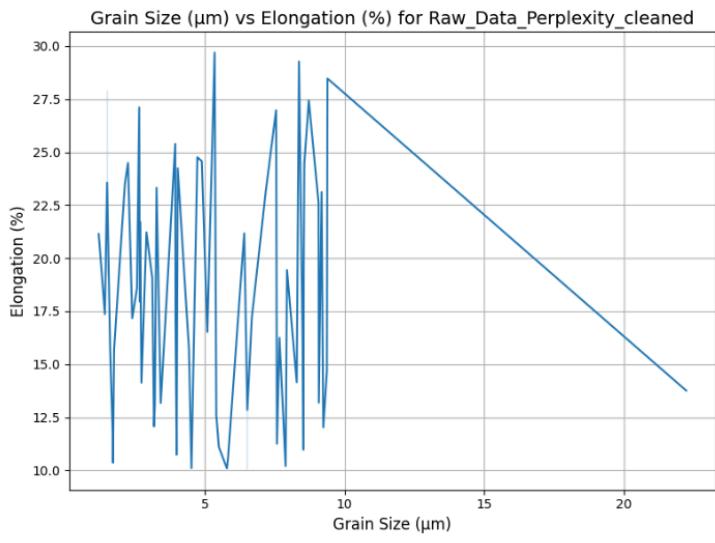
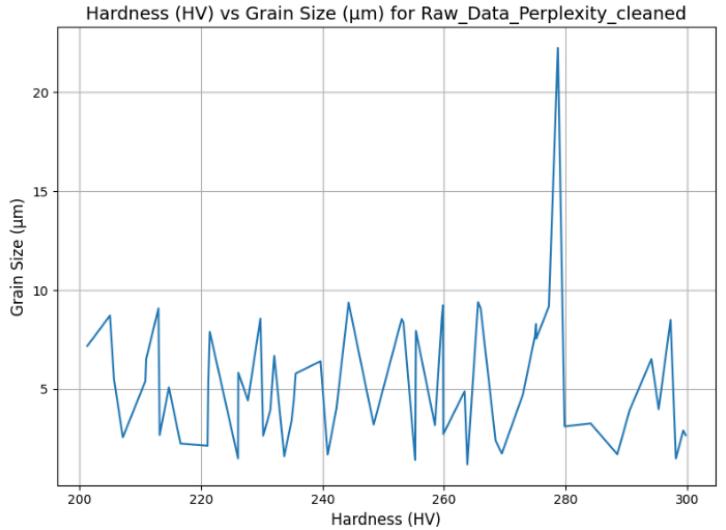
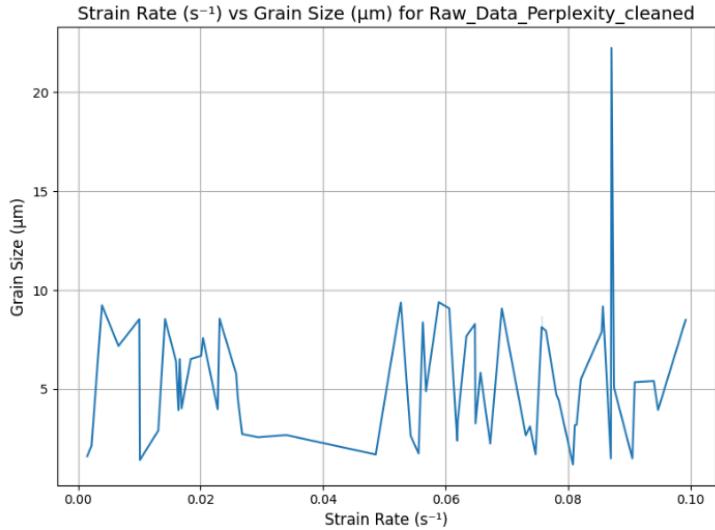
Pair Plot of Key Properties for Raw_Data_Gemini_cleaned.csv





Perplexity: High variance in dislocation density and mechanical properties. Inconsistencies more visible here but still within physically plausible range.





Conclusion

The AI-generated datasets largely exhibited physically meaningful values, adhering to the known trends and behaviors of duplex steels—particularly in the correlation between microstructure and mechanical strength. After cleaning, no unphysical values (such as negative hardness or temperature) were present.

The datasets offered broad yet realistic value ranges for key mechanical properties such as yield strength, ultimate tensile strength, hardness, grain size, and strain rate. These ranges are suitable for simulating varied processing conditions in duplex stainless steels, allowing flexibility and generalization.

Inconsistencies within the data, such as unit variations, missing values, and typographical errors, were generally mild to moderate. Most were syntactic in nature—like including units (e.g., “MPa”) within numeric fields—and

were effectively resolved through systematic data cleaning. These issues also provided a practical test of preprocessing workflows.

Notable property relationships were captured across most datasets:

- **Strain rate vs. grain size** illustrated the impact of deformation on microstructure.
- **Grain size vs. hardness** reflected an inverse correlation due to the Hall-Petch effect.
- **Grain size vs. elongation** showed how coarser grains improve ductility at the expense of strength.
- **Grain size vs. yield strength** confirmed the trend of increased strength with grain refinement.

These consistent trends across multiple AI platforms suggest strong potential for using synthetic datasets to replicate realistic mechanical behavior.

Overall, AI tools like ChatGPT, Claude, Gemini, and Perplexity have demonstrated the ability to generate high-quality synthetic data for materials science applications. When properly cleaned and validated, this data can be effectively used for machine learning model development, exploratory analysis, and accelerating materials research. Despite the need for human oversight to filter out occasional formatting or logic errors, AI-generated data holds significant promise in supplementing experimental datasets and enabling rapid, data-driven discovery.

Comparison with Original Dataset

