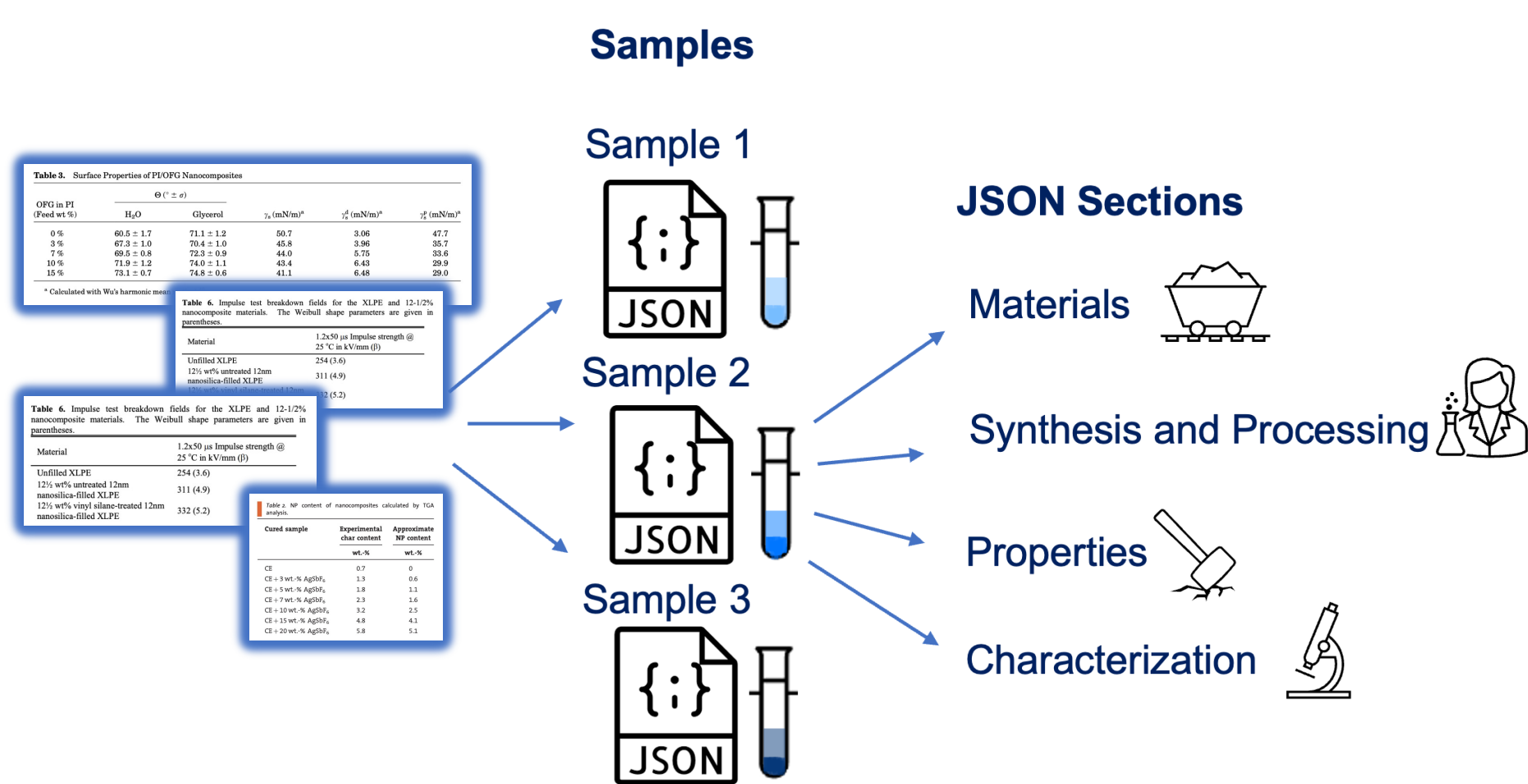# Using Large Language Models for Data Extraction from Tables in Materials Literature

**Defne Circi**, Ghazal Khalighinejad, Anlan Chen

Bhuwan Dhingra, L. Catherine Brinson

Duke UNIVERSITY

## Unlocking Insights in Scientific Literature

- **Finding and using data from literature** is a common problem.
- We need to **search among many documents** for key information.
- Traditionally, data extraction is done **manually → time consuming** and **tedious**
- Collecting experimental data **at a scale** is critical.
- **Large language models (LLMs)** can make the information most important to scientists, such as material identification and properties readily available.
- Composition and properties of materials are predominantly **condensed in tables**.

**Unstructured**  **Structured**



## 1. Goal: Extracting multiple experimental samples per table



**Samples**

Sample 1 · Sample 2 · Sample 3

**JSON Sections**
- Materials
- Synthesis and Processing
- Properties
- Characterization

## 2. Dataset overview and ground truth samples

- Articles from **MaterialsMine** database
- **Table dataset**: 18 articles, 37 tables and their captions, 182 samples
- Sample size range from 2 to 15
- On average 3.1 material properties in each table

## 3. Choosing inputs of table data

### Option 1: GPT-4-Vision on table image



**Text + IMAGE**  GPT-4(V)

### Option 2: GPT-4 on unstructured OCR extraction from table image



**Text + OCR**  GPT-4

### Option 3: GPT-4 on unstructured OCR extraction from table image



**Text + STRUCTURED FORMAT**  GPT-4

## 4. Evaluation of LLM output

### Composition level evaluation



**Ground Truth**

Sample id: 1,
matrix name: PP,
filler name: silica,
composition: {amount: 5%, type: wt},
particle surface treatment name: not specified,

**Predictions**

Sample id: 1,
matrix name: PP,
filler name: none,
composition: {amount: 0.0%, type: wt},
particle surface treatment name": not specified,

match — not a match — partial match — match

### Accuracy scores of composition information extraction

| Input type/Including missing samples | no | yes |
|---|---|---|
| Image | $0.917 \pm 0.036$ | $0.910 \pm 0.037$ |
| OCR | $0.890 \pm 0.065$ | $0.790 \pm 0.107$ |
| Structured Format (with captions) | $0.948 \pm 0.032$ | $0.816 \pm 0.113$ |
| Structured Format (without captions) | $0.890 \pm 0.056$ | $0.832 \pm 0.089$ |

### Property and condition level evaluation



**Ground Truth**

properties: {
  Example Property Identifier: {value: 910,
  unit: MPa, conditions: [{type: temperature, value: -413, unit: K}],
  Another Key Name: {...}
}

**Predictions**

properties: {
  Not Close Property Name: {...},
  Example Property Name : {value: 910, unit: MPa, conditions: [{type: temperature, value: -413, unit: K}],
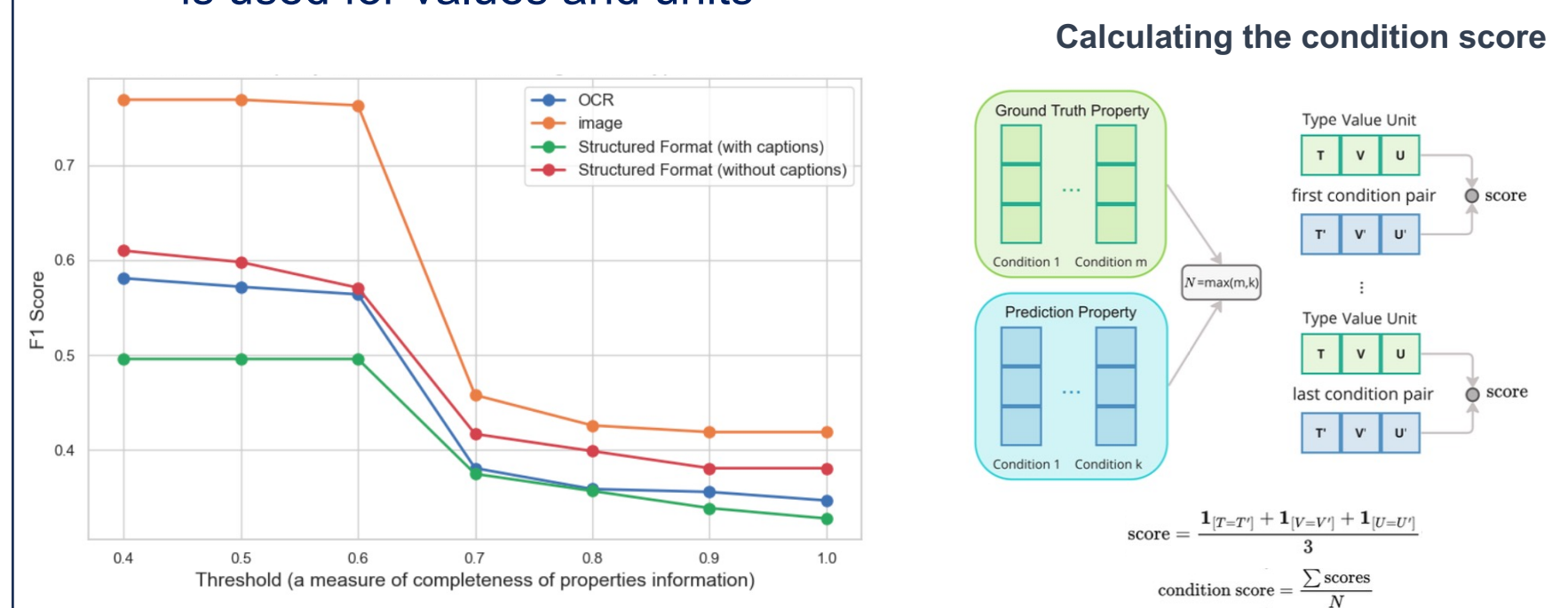  }

property match

### F1 scores of property name information extraction

| Input type/Including missing samples | no | yes |
|---|---|---|
| Image | $0.869 \pm 0.078$ | $0.863 \pm 0.078$ |
| OCR | $0.766 \pm 0.104$ | $0.666 \pm 0.117$ |
| Structured Format (with captions) | $0.795 \pm 0.107$ | $0.682 \pm 0.129$ |
| Structured Format (without captions) | $0.617 \pm 0.133$ | $0.576 \pm 0.134$ |

### F1 scores of property information considering value, unit and condition
- Calculated a matching score for each of the entities. The final score for a property is an average of these individual scores. Equality check is used for values and units



**Calculating the condition score**

$$score = \frac{1_{[T=T']} + 1_{[V=V']} + 1_{[U=U']}}{3}$$

$$\text{condition score} = \frac{\sum scores}{N}$$

## Findings

- Multimodal model with an image output yielded the most promising results.
- We introduced a flexible evaluation technique tailored to assess the accuracy and efficiency of these extraction methods, contributing to a nuanced understanding of their performance on this complex task.

## Future work

- **Granular benchmarking** across entity and relationship types
- **Benchmarking** across commercial and open-source models
- Extracting sample information from **tables, figures and text**
- Scaling complex extraction and verification to **various materials domains**