# Extracting accurate data from research papers

Using Large Language Models (LLMs)

**Maciej P. Polak,** Dane Morgan

Department of Materials Science and Engineering
University of Wisconsin – Madison

June 20th 2024

# Our approach to data extraction

❖ Break up papers into sentences and **classify**
  ○ *Does it contain relevant data?*

❖ **Extract** the data
  ○ By hand
  ○ Automatically – with **LLMs**

❖ **Accurate** extraction is possible:
  ○ ~90% precision and recall

# Our approach to data extraction

❖ **Problem:** high quality data needed for building machine learning models
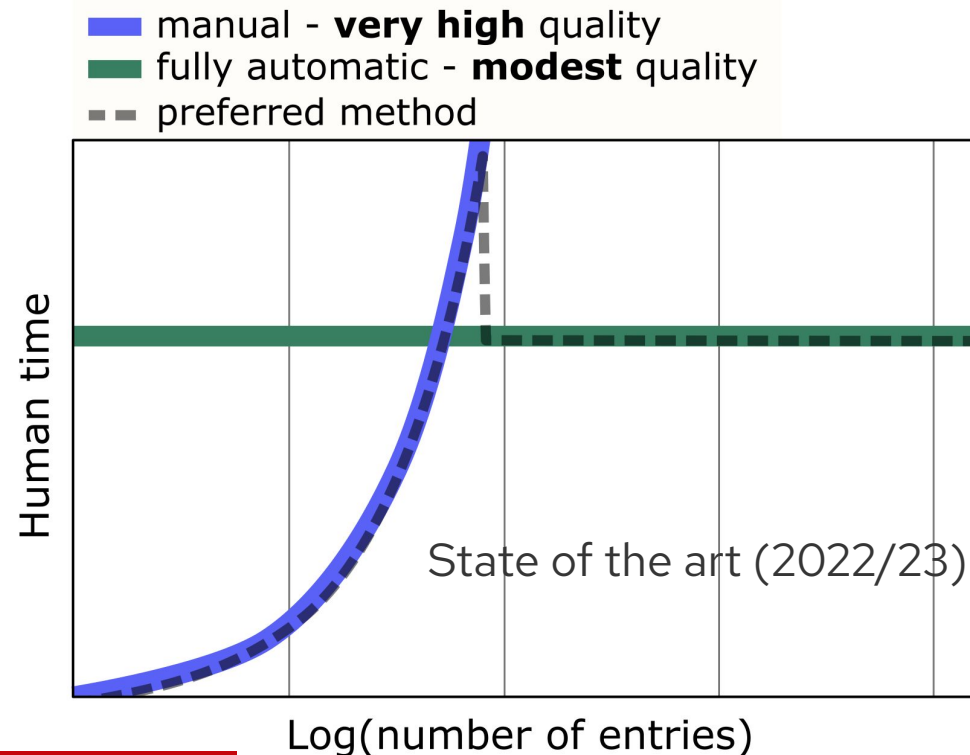❖ **Solution:** extract data from research papers
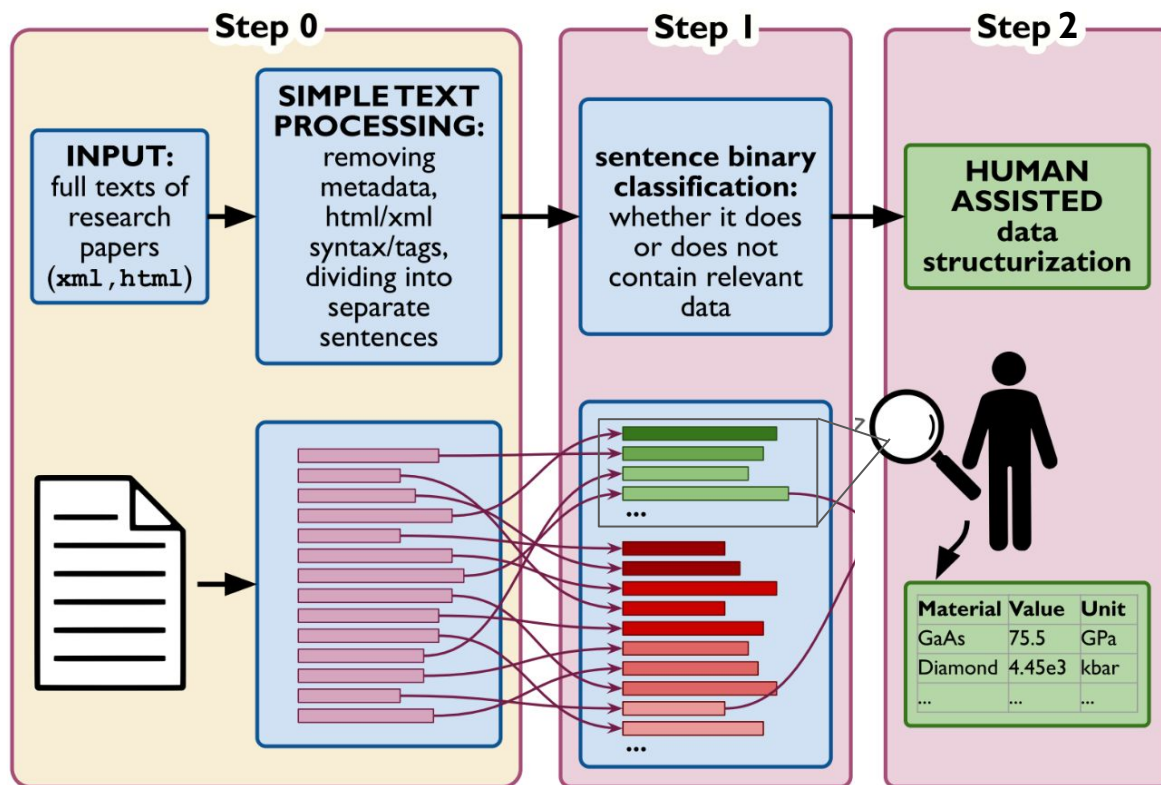- ○ **Problems:**
  - ■ low quality of auto–mated data extraction
  - ■ too many papers to extract manually
- ○ **Solution:**
  - ■ *use language models to extract data from research papers*



Legend:
- manual - **very high** quality
- fully automatic - **modest** quality
- preferred method

State of the art (2022/23)

Human time

Log(number of entries)
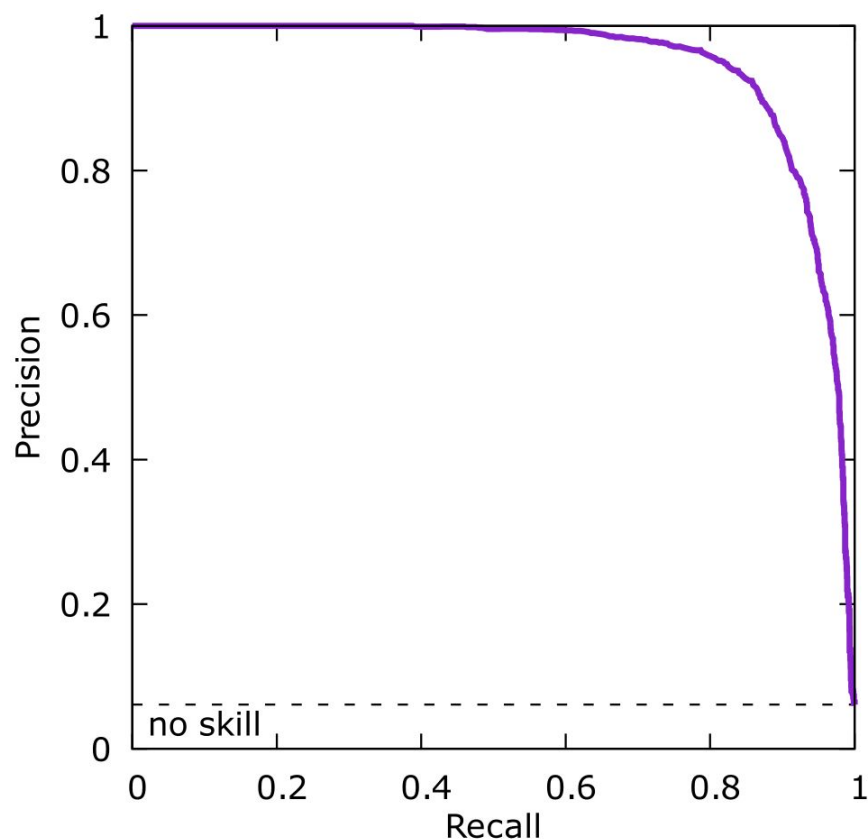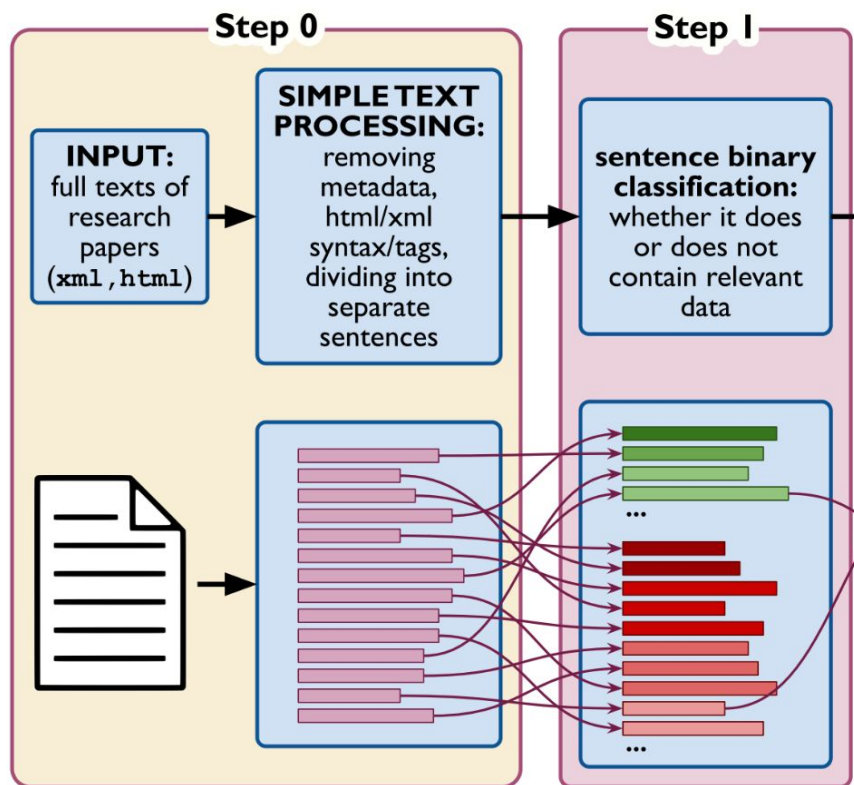
# Using (L)LMs for sentence classification



**Data triplet:**
Material, Value, Unit

# Using (L)LMs for sentence classification

# Flexible, model-agnostic method for materials data extraction from text using general purpose language models

Maciej P. Polak, [iD] * Shrey Modi, [iD] Anna Latosinska, Jinming Zhang, Ching-Wen Wang, Shaonan Wang, Ayan Deep Hazra and Dane Morgan*

❖ Extremely simple
❖ Needs minimal resources
❖ Almost no coding required
❖ Used to develop the most complete and largest to date database of critical cooling rates of metallic glasses



Legend:
- manual - **very high** quality
- this work - **high** quality
- fully automatic - **modest** quality
- preferred method

Axis labels: Human time (y-axis), Log(number of entries) (x-axis)

5

# Automate data structurization with LLMs



Using **GPT-4**:
- ~90% recall
- ~30% precision

# Automate data structurization with LLMs



Using **GPT-4**:
- ~98% recall
- ~42% precision

# Automate data structurization with LLMs



Using **GPT-4**:
- ~88% recall
- ~91% precision
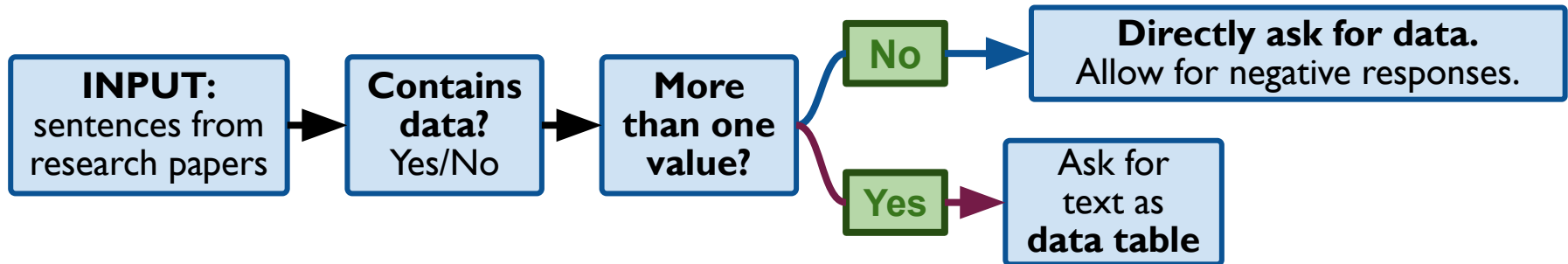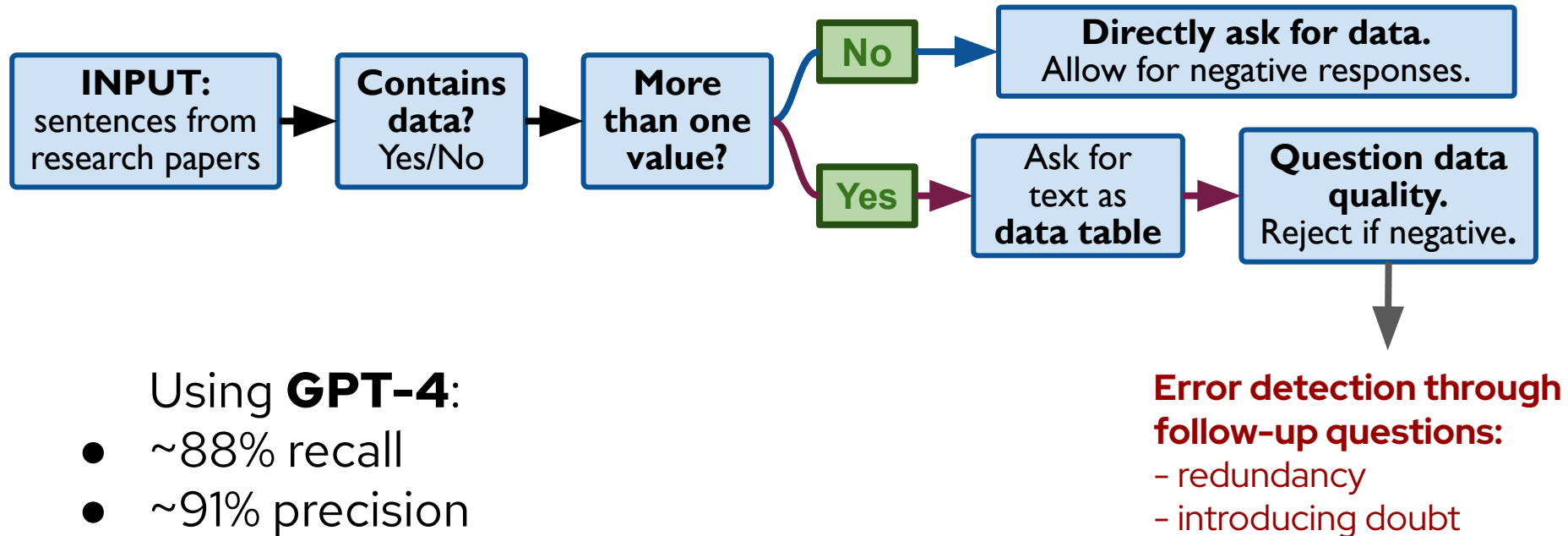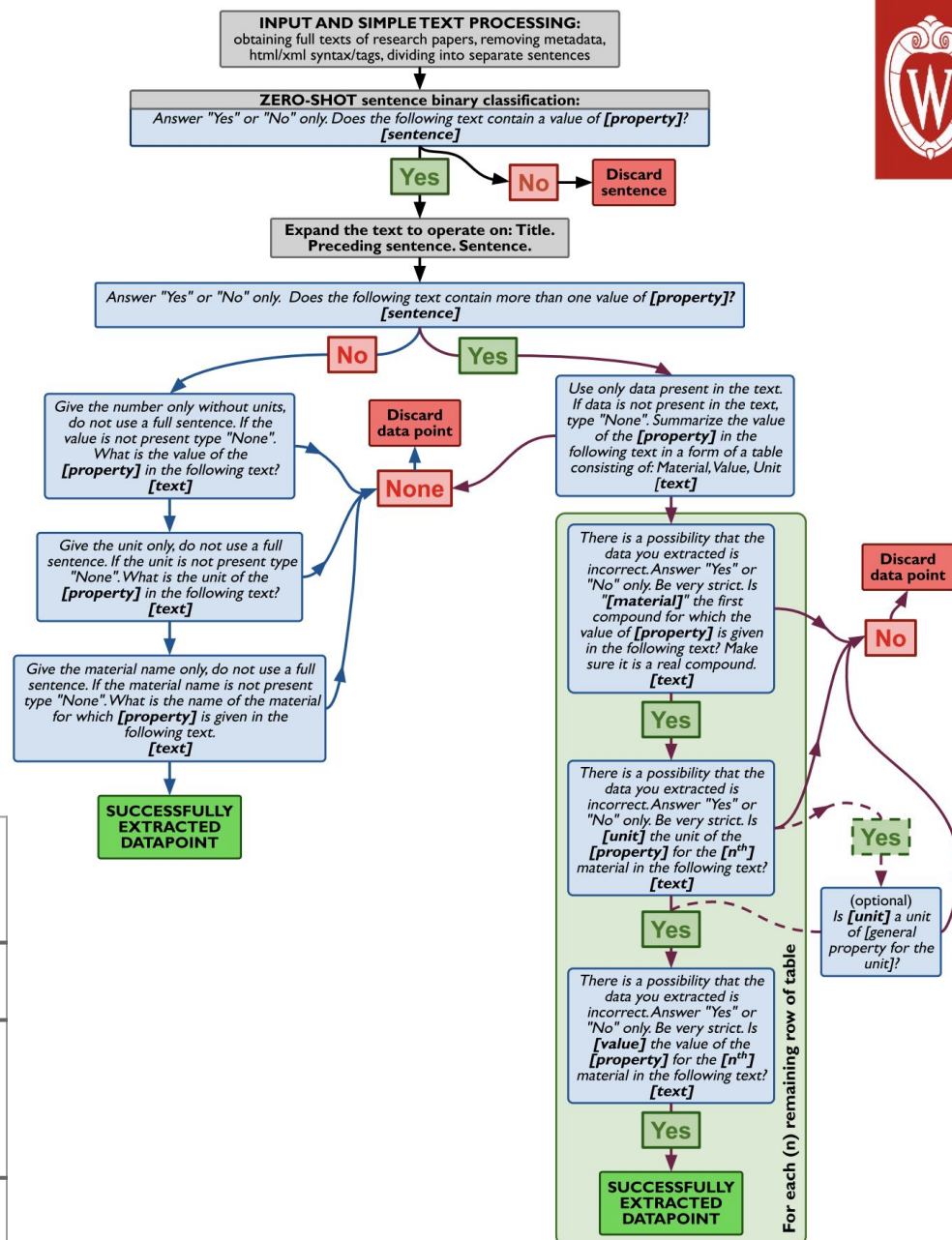
# ChatExtract

a workflow consisting of a series of engineered conversational prompts to a large language model (LLM), and actions based on the model's responses, with error detection.

| | Precision (%) | Recall (%) |
|---|---|---|
| ChatExtract (GPT4) | **90.8** | **87.7** |
| Chain-of-thought (GPT4) (ChatExtract without error-detection) | 42.7 | 98.9 |
| Previous non LLMs | <50% | <50% |



**INPUT AND SIMPLE TEXT PROCESSING:** obtaining full texts of research papers, removing metadata, html/xml syntax/tags, dividing into separate sentences

**ZERO-SHOT sentence binary classification:** Answer "Yes" or "No" only. Does the following text contain a value of *[property]*? *[sentence]*

Yes → **Expand the text to operate on: Title. Preceding sentence. Sentence.**
No → **Discard sentence**

Answer "Yes" or "No" only. Does the following text contain more than one value of *[property]*? *[sentence]*

No / Yes

Give the number only without units, do not use a full sentence. If the value is not present type "None". What is the value of the *[property]* in the following text? *[text]*

Give the unit only, do not use a full sentence. If the unit is not present type "None". What is the unit of the *[property]* in the following text? *[text]*

Give the material name only, do not use a full sentence. If the material name is not present type "None". What is the name of the material for which *[property]* is given in the following text. *[text]*

**None** / **Discard data point**

**SUCCESSFULLY EXTRACTED DATAPOINT**

Use only data present in the text. If data is not present in the text, type "None". Summarize the value of the *[property]* in the following text in a form of a table consisting of: Material, Value, Unit *[text]*

There is a possibility that the data you extracted is incorrect. Answer "Yes" or "No" only. Be very strict. Is *"[material]"* the first compound for which the value of *[property]* is given in the following text? Make sure it is a real compound. *[text]*

Yes

There is a possibility that the data you extracted is incorrect. Answer "Yes" or "No" only. Be very strict. Is *[unit]* the unit of the *[property]* for the *[n^{th}]* material in the following text? *[text]*

Yes

There is a possibility that the data you extracted is incorrect. Answer "Yes" or "No" only. Be very strict. Is *[value]* the value of the *[property]* for the *[n^{th}]* material in the following text? *[text]*

Yes

**Discard data point**

No

**Yes** (optional) Is *[unit]* a unit of [general property for the unit]?

For each (n) remaining row of table
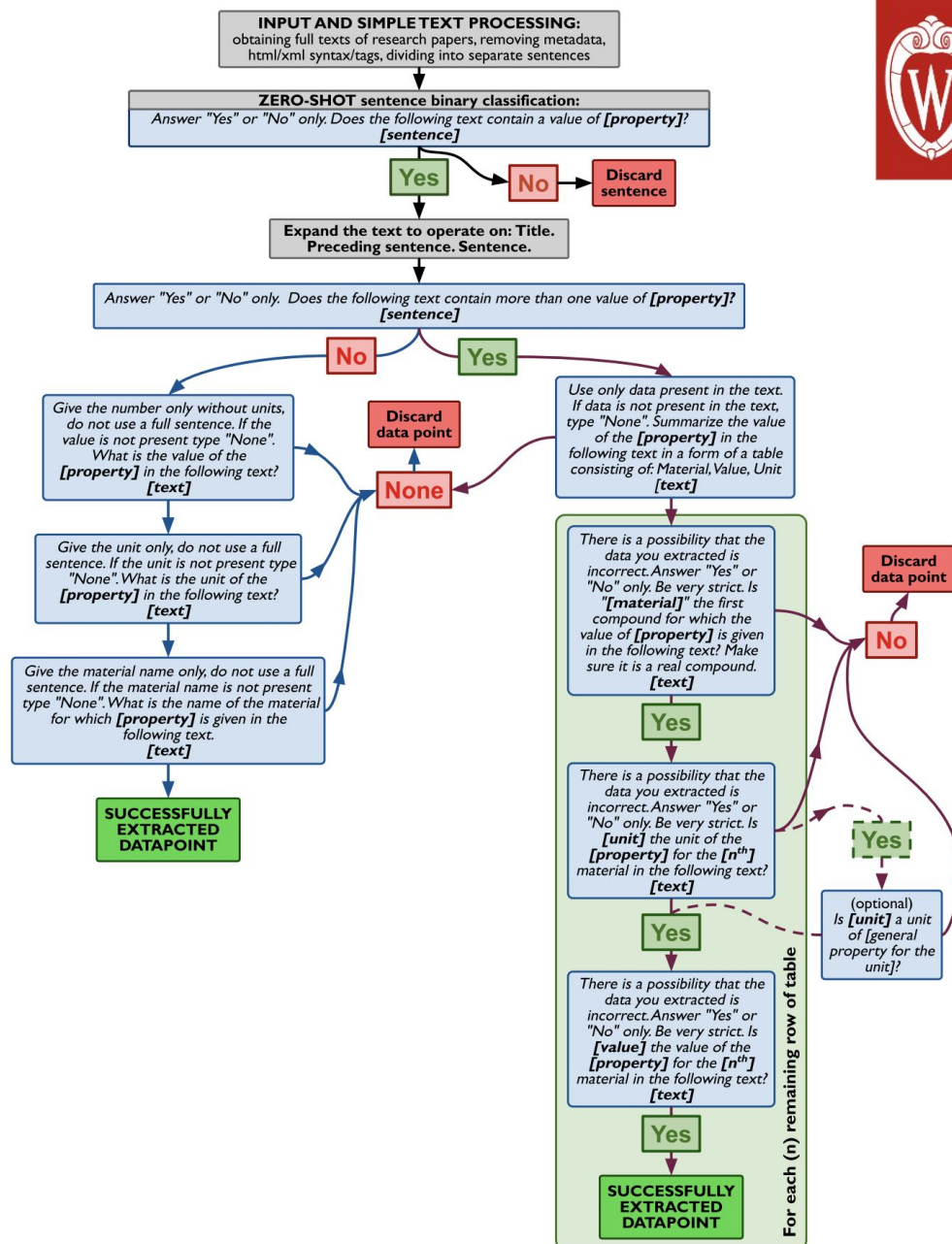
**SUCCESSFULLY EXTRACTED DATAPOINT**

# ChatExtract

**Key concepts:**

- Break up tasks into the simplest possible steps
- Strictly enforce output format
- Allow for negative answers
- When validating – be redundant, introduce doubt

**Benefits of our approach:**

- Flexible and adaptable
- No need for deep understanding of data to be extracted

# Extracting accurate materials data from research papers with conversational language models and prompt engineering

Maciej P. Polak [1] ✉ & Dane Morgan [1] ✉

## Metallic Glass Critical Cooling Rates

- Returned 684 papers, 110,126 sentences.
- ChatExtract (GPT-4) extracted 721 values.
- Standardized version had 280 values (120 unique compounds), 1.5x previous databases.

## High-Entropy Alloy Yield Stress

- 4029 research papers, 840,431 sentences.
- ChatExtract extracted 8,961 values.
- Standardized version had 2416 values.
- Largest database up to date

## Performance

- Precision and Recall ~90%
- ~10m human time
- ~3h compute time (ChatGPT)
- ~0-2h Standardization



High-Entropy Alloys

Number of Alloys vs Range of Yield Stress Values (MPa)