# Language Models Enable Simple Systems for Generating Structured Views of Heterogeneous Data Lakes

Simran Arora
Stanford University
simarora@stanford.edu

Brandon Yang
Stanford University
bcyang@stanford.edu*

Sabri Eyuboglu
Stanford University
eyuboglu@stanford.edu*

Avanika Narayan
Stanford University
avanikan@stanford.edu

Andrew Hojel
Stanford University
ahojel@stanford.edu

Immanuel Trummer
Cornell University
itrummer@cornell.edu

Christopher Ré
Stanford University
chrismre@stanford.edu

## ABSTRACT

A long standing goal in the data management community is developing systems that input documents and output queryable tables without user effort. Given the sheer variety of potential documents, state-of-the art systems make simplifying assumptions and use domain specific training. In this work, we ask whether we can maintain generality by using the in-context learning abilities of large language models (LLMs). We propose and evaluate Evaporate, a prototype system powered by LLMs. We identify two strategies for implementing this system: prompt the LLM to directly extract values from documents or prompt the LLM to synthesize code that performs the extraction. Our evaluations show a cost-quality tradeoff between these two approaches. Code synthesis is cheap, but far less accurate than directly processing each document with the LLM. To improve quality while maintaining low cost, we propose an extended implementation, Evaporate-Code+, which achieves better quality than direct extraction. Our insight is to generate many candidate functions and ensemble their extractions using weak supervision. Evaporate-Code+ outperforms the state-of-the art systems using a *sublinear* pass over the documents with the LLM. This equates to a 110× reduction in the number of documents the LLM needs to process across our 16 real-world evaluation settings.

## 1 INTRODUCTION

Organizations often seek insights trapped in heterogeneous data lakes (*e.g.* the web, corporate data lakes, and electronic health records) [10, 26, 58]. In their raw form, these data sources cannot

easily support analytical queries. A long standing goal of the data management community is to develop systems that automatically convert heterogeneous data lakes into queryable, structured tables [12, 15, 50, 70, inter alia.]. In this work, we investigate whether recent large language models can help address this problem.

We study systems that take as **input** heterogeneous documents (*e.g.* HTML webpages, PDFs, text) and **output** a tabular, structured view of the documents. These systems must identify the schema and perform extraction to populate the table.

> EXAMPLE 1. Medical researchers frequently use data spanning electronic health records (EHR), clinical trials, knowledge sources (e.g. PubMed), and FDA reports to understand and monitor patients and treatments [8]. Consider the large collection of **FDA 510(k)** reviews for premarket medical devices, which have been the subject of multiple studies [68, 72]. Our objective is to output a table that automatically structures the attributes that are distributed in the ~20-page **PDFs**, for instance the `device classification`, `predicate device code`, and `indications for use`.

Systems designed to tackle this problem must balance a three-way tradeoff between **cost** (data lakes may hold millions of documents), **quality** (output tables should be able to accurately support an analyst's queries), and **generality** (different data lakes have different document types and structure). See Section 2 for a formal task definition and further discussion of this tradeoff.

Given the range of formats, attributes, and domains across documents, prior systems rely on simplifying assumptions (e.g. handling one document format). The majority of works focus on structuring HTML [12, 15, 25], assuming the attributes and values are at specific positions in the HTML-DOM [22, 45, 46, 71]. For unstructured text, current approaches use linguistic tools (e.g., dependency parsers) to introduce structure [15, 25, 31, 51] and then apply heuristic rules over the resulting structure to extract information. The documents in Example 1 highlight the limitations of the prior approaches: they lack (e.g. HTML) structure and, consistent with recent evaluation efforts [71], we find the SoTA approaches for unstructured text perform poorly on long semi-structured PDFs (See [5]). Some systems assume there is a human-in-the-loop, labeling data and writing heuristic rules for extraction [57, 61], while others assume access to annotated training documents from the domain [22, 45, 46]. Researchers manually annotated the reports in Example 1 [68].
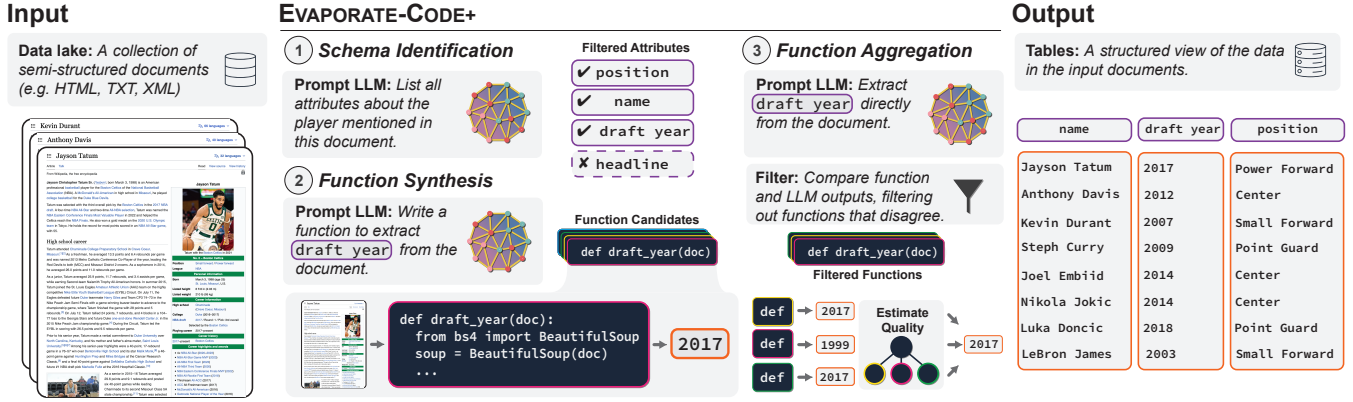
**Figure 1: The user provides a collection of documents (e.g. NBA player bios) and EVAPORATE outputs a table by identifying attributes and populating columns. EVAPORATE avoids running expensive LLM inference on all documents by (1) synthesizing the key attributes from a small sample of documents and (2) synthesizing (e.g. Pythonic) functions that then are reused at scale to process documents. Because function quality is variable, EVAPORATE (3) applies an algorithm that generates many candidate functions and ensembles their extractions using weak supervision.**

In this work, we explore whether we can improve generality by leveraging *large language models* (LLMs). An LLM is a deep learning model that is pretrained on broad data and can be adapted to diverse tasks, from machine translation to data wrangling [14, 49]. At inference time, the models take as input a natural language task description termed a *prompt* [11, 14] and generate a natural language response. See Section 2.3 for more background on LLMs.

**EVAPORATE.** (Section 3) We present EVAPORATE, a system that uses LLMs to produce structured views of semi-structured data lakes. Our evaluation spans 16 real-world settings from movie and university websites to *FDA 510(k)* reviews [22, 32, 34, 37, 45, 68, 72].

The user inputs a collection of documents and EVAPORATE automatically identifies the schema and performs extraction to populate the table. Our implementation requires *no customization, training, or human effort* to support the diverse evaluation settings. We propose two fundamental strategies for implementing this interface, identifying a tradeoff between their cost and quality:

(1) EVAPORATE-DIRECT (Figure 2) The LLM directly extracts values from documents.
(2) EVAPORATE-CODE (Figure 4) The LLM *synthesizes code* that is then applied to process documents at scale.

EVAPORATE-CODE is cheap, but underperforms EVAPORATE-DIRECT by 24.9% (13.8 F1 points) averaged across our evaluation settings. We thus seek a new code synthesis approach. We present EVAPORATE-CODE+, which achieves better quality than direct extraction. Our insight is to synthesize many code snippets for extraction and ensemble their outputs using weak supervision.

**Direct Extraction (Section 3.1).** Our first implementation, EVAPORATE-DIRECT, applies a *single prompt* (included in [5]) to each document in the input. The prompt instructs the LLM to both identify the schema and extract values. Remarkably, we find that in some settings, with a single prompt and no task specific modifications, performance is already competitive with state-of-the-art systems that rely on domain specific assumptions and training.

However, this implementation is very expensive. LLMs are optimized for interactive, human-in-the-loop applications (*e.g.* ChatGPT) [69], not high-throughput data processing tasks [60]. The number of tokens processed by an LLM in EVAPORATE-DIRECT grows *linearly* with the size of the data lake. As of March 2023, applying OpenAI's models to the 55 million Wikipedia articles would cost over $110k (gpt-3.5, $0.002/1k tokens) and $1.1M (text-davinci-003, $0.02/1k tokens) dollars [1, 52]. There are *billions* of webpages on the broader Internet [35] and the facts change over time. For instance, NBA players are added to Wikipedia, a player's team changes after trades, and the points per game metric changes after every game. Data processing is a *routine expense* (repeated by multiple data analysts), not a one-time cost [59].

**Code Synthesis (Section 3.2).** *Can we produce the structured table using a sublinear pass of the LLM over the documents?* We propose EVAPORATE-CODE, which splits the task into two sub-tasks: (1) identify the table schema and (2) extract values. This view allows us to exploit the distinct *redundancies* of each sub-task that occur when running LLM inference on every document:

(1) *Schema Generation.* In order to identify a schema, we only process a small sample of documents with the LLM. This succeeds because there is redundancy in the attributes mentioned across documents. For e.g., in Example 1, most reports mention a predicate device name.
(2) *Function Synthesis.* We prompt the LLM to synthesize (*e.g.* Pythonic) *functions*, that can be applied at scale across the documents. This works because of redundancy in the formatting of attribute-value pairs. For e.g., the FDA 510(k)s use the consistent format "Predicate device name: k".

The number of tokens processed by the LLM in EVAPORATE-CODE is *fixed* and does not grow with the size of the data lake (as illustrated in Figure 3), addressing the cost issues of EVAPORATE-DIRECT. However, the LLM synthesizes variable quality information extraction functions. The extractions are up to 14 points worse in Pair F1 score than those produced using EVAPORATE-DIRECT.

**Code Synthesis + Aggregation.** (Section 3.3) To improve quality while keeping costs low, we propose EVAPORATE-CODE+. Studying the synthesized functions, we observe some only work for a narrow slice of documents, while others exhibit syntactic and logical errors. To reduce variance, we synthesize many candidate functions, then estimate their quality and aggregate their extractions using *weak supervision*. This builds on our work [4], which broadly applies weak supervision to prompting for the first time.

Weak supervision (WS) is a statistical framework for modeling and combining noisy sources with varied coverages without any labeled data [57, 66]. However, WS is typically applied over *human-generated* functions while our setting consists of *machine-generated* functions. This presents issues when attempting to apply existing WS tools. (1) WS theoretically assumes all noisy sources are better than random performance (50% accuracy), yet 40% of our generated functions are *below 25%* (Section 3.2). (2) WS attempts to deploy functions that achieve high quality on narrow slices of data (high precision), and allow the function to *abstain* on data external to the slice (low recall). While humans can express when functions should abstain, the machine-generated functions do not contain this logic. To handle the *open* WS setting, we introduce a new algorithm for ensembling the functions (Algorithm 1).

We summarize our overall contributions as follows.

(1) **Our system offers new capabilities for the long-studied structured view generation problem.** Existing systems require in-domain training and handle limited document formats (e.g. HTML [16, 22, 24, 45]). EVAPORATE requires no training and succeeds on different document formats (HTML, PDF, TXT) off-the-shelf. (Section 4).

(2) **We study a new tradeoff space between direct extraction and code synthesis for data tasks.** EVAPORATE *asymptotically* reduces the number of tokens that the LLM needs to process to generate the outputs. At 10k documents per evaluation setting, this amounts to a 110x cost reduction. Further, prior works using LLMs for data tasks require users to manually write prompts [49]. EVAPORATE is built with task-agnostic prompts that generalize across settings.

(3) **We present an algorithm and theoretical analysis for applying weak supervision to open-ended functions and extraction tasks.** Although EVAPORATE-CODE is more efficient, EVAPORATE-DIRECT achieves significantly higher quality. Using our algorithm, EVAPORATE-CODE+ outperforms EVAPORATE-DIRECT, which directly processes every document, by 10.1 F1 points (18%) (Table 3).

(4) **We extensively validate the system on 16 data settings from 5 domains and 3 data formats, and across 4 LLMs.** (1) EVAPORATE outperforms the SoTA *learned* baseline systems by 3.2 F1 points (6%) when generating tables (both schema generation and extraction) end-to-end, and 6.7 F1 points (10%) on the extraction step. (2) EVAPORATE-CODE+ achieves a 10.1 F1 point increase over EVAPORATE-DIRECT, using text-davinci-003. (3) Across four unique LLMs we show the relative quality of EVAPORATE-DIRECT vs. EVAPORATE-CODE+ remains consistent.

We define the problem in Section 2. We present EVAPORATE in Section 3, evaluations in Section 4, and related works in Section 5.

## 2 PRELIMINARIES

We first define the problem setting and system desiderata.

### 2.1 Problem Setting

We study the problem of constructing a structured view (*i.e.* database table) of a set of semi-structured documents (*e.g.* HTML, PDF, TXT). Formally, we define the problem as follows:

- **Input:** User provides a set of $n$ semi-structured documents $D = \{d_1, d_2, ...d_n\}$ (*e.g.* A collection of FDA 510(k) reviews for premarket notification submission for medical devices).
- **Output:** System outputs a table defined by a set of attribute names $A = \{a_1, a_2, ...a_m\}$ (*e.g.* $a_1$ =indications for use, $a_2$=classification) and a set of $n$ extracted records for $R = \{r_1, r_2, ...r_n\}$, one per document, where $r_i$ is an $m$-tuple (*e.g.* $r_1$ = ("fracture", "x-ray")).

Unlike prior work which proposes systems that rely on manual labeling [61] or manual prompt tuning [49, 63], we aim to develop *automated* solutions, which require no user effort.

*Measuring System Quality* We compare the generated table $(A, R)$ to a manually curated "ground-truth" table $(\hat{A}, \hat{R})$. The coverage of an attribute refers to the fraction of documents that include the attribute and its value. Following prior work, we prioritize attributes with high *coverage*, which tend to be useful for analysis [16, 18]. We measure agreement between the tables using Pair F1. For additional details on our evaluation setup, see Section 4.3.

### 2.2 System Desiderata

Current systems for producing structured views are limited in their generality, cost/flexibility, and quality/usability [16, 18, 51, 71]. Here we review the existing systems.

**Generality.** *The ideal system will generalize across document formats and domains, without manually engineered rules or task-specific training.* This is important because the input documents $D$ could focus on any imaginable topic or use any file format [71]. Existing systems featurize documents by tagging the named entities (NER), dependency parse tree, and part-of-speech (POS), and train a model to predict whether a span of text is a useful fact [39]. Unfortunately, the performance of the parse, NER, and POS tags drastically degrade on semi-structured data (e.g. HTML elements) and longer sequences of text (i.e. full documents) [71]. We provide detailed error analysis [5]. A specialized class of systems focuses on processing semi-structured web HTML documents by leveraging the HTML DOM tree as features [13, 16, 22, 25, 46, inter alia.]. However, the systems thus do not support other document formats.

**Cost.** *The ideal system will enable users to manage a cost-coverage tradeoff, rather than requiring them to extract "all-or-nothing".* The existing systems are built to extract *all* possible facts in the documents, without prioritizing important attributes or allowing the user to influence what is extracted [21, 71]. Processing every line of every document can be expensive. To mitigate this, the user can define the attributes of interest then apply a closed IE system for extraction, however this requires upfront human effort. **Desiderata**: The ideal system will enable users to manage a cost-coverage tradeoff, rather than requiring them extract "all-or-nothing".

**Quality.** *The ideal system will output a table* $(A, R)$ *with full columns (i.e. high-coverage attributes) and accurate, consistently formatted extractions.* Existing OpenIE systems commonly extract tuples in unnormalized forms directly from documents [21]. This can make the resulting extractions difficult to use for analysis, requiring advanced systems or user-defined post-processing code for resolving subject, objects, and predicates to a canonical form [15].

## 2.3 Background on Large Language Models

In this section, we provide background on *large language models* (LLMs), which are central to our work.

DEFINITION 1 (LARGE LANGUAGE MODEL). *A machine learning model,* $\mathcal{F}$*, trained on a self-supervised task (e.g. next word prediction) over a massive corpus of text [29]. Language models can be used to generate new text based on provided context. For example:*

$$\mathcal{F}(\textit{All that glitters}) \rightarrow \textit{is not gold}.$$

Numerous studies have demonstrated LLMs capable of solving new tasks without updating any model parameters, a phenomenon termed *in-context learning* [2, 14, 49]. Specifically, these studies show that when passed an appropriate description of the task, the model often generates text completing the task.

DEFINITION 2 (PROMPT). *A natural language task-specification used to elicit a particular generation from an LLM. Prompts often include demonstrations of the task. For example, the prompt below elicits the translation of the word cheese into French:*

$$\underbrace{\mathcal{F}(\textit{Translate. Eng: hello, Fr: bonjour; Eng: cheese, Fr: })}_{\textit{Prompt}} \rightarrow \underbrace{\textit{fromage}}_{\textit{Generation}}$$

Examples of prompts used in this work are provided in Figures 2 and 4. All prompts used in the system are provided in [5].

## 3 EVAPORATE: A PROTOTYPE SYSTEM POWERED BY LANGUAGE MODELS

We introduce EVAPORATE, a prototype system that uses LLMs to materialize a structured view of a heterogeneous, semi-structured data lake. Compared to prior systems, which rely on manual labeling [61] or tuning prompts to a domain [49], EVAPORATE exposes a remarkably *general* interface: the user inputs documents and the system automatically outputs a structured view of those documents, without any domain specific training or prompt customization.

*Overview.* We instantiate the EVAPORATE interface with three different implementations. We can feed every document to the LLM and prompt it to extract values directly (*direct extraction*, Figure 2), or feed a small sample of documents to the LLM and prompt it to write *code* to do the extraction (*code extraction*, Figure 4). In Section 3.1 and Section 3.2, we describe baseline implementations of these two strategies, EVAPORATE-DIRECT and EVAPORATE-CODE. We find that these two implementations tradeoff cost and quality. Then, in Section 3.3, we propose a code extraction implementation that uses weak supervision to improve quality and retain low cost.

*Prompt Management* EVAPORATE applies a set of task-agnostic prompts, which are all provided verbatim in the technical report [5]. These tasks are not modified for different tasks. Within the system, the prompts are Python f-strings, with placeholders for



**Figure 2: Prompt for EVAPORATE-DIRECT structured. The prompt template, which includes placeholders for in-context examples and and the inference example (i.e., data lake documents), is applied to each document in the data lake.**

inputs chunks of text from the particular dataset being processed. The LLM is prompted with the formatted strings. We use a caching tool that we helped develop called MANIFEST [53] to store input and completion pairs from the LLM prompting in a local SQLite database, where keys are the prompt-inputs and values are the completions . Therefore, if users repeatedly run the system on the same dataset, they do not incur the LLM inference costs again.

## 3.1 EVAPORATE-DIRECT

In this section, we describe a simple *direct extraction* implementation, EVAPORATE-DIRECT that applies a single prompt template to every document. This prompt template, which is included in Figure 2, instructs the LLM to both identify the schema and extract values (see [5] for the full prompt). It consists of a few in-context examples that are general, *i.e.* are not customized to a particular format, domain, or document.

Below we discuss how we (1) manage long documents that cannot fit in the LLM's context window, (2) process the LLM's textual outputs, (3) prioritize the most useful attributes according to principles described in prior work [16].

*Managing long documents.* The input to EVAPORATE is a file path to raw documents, which can be several pages long. For instance the Medical FDA reports in Example 1 are ~20 pages long. However, the underlying Transformer architecture of modern LLMs is limited to processing a fixed number of tokens (*e.g.* a few thousand tokens), referred to as the *context window*, during each inference call. EVAPORATE therefore splits the raw documents such that each piece is within the context window. Each chunk is inserted into the prompt in turn as shown in Figure 2.

*Processing text outputs.* Language models output open ended text so the last step is to convert this to a usable table. To facilitate this data transformation, we can specify formats in our prompt demonstrations to encourage the LLM to organize the output in a similar structure. For instance, the demonstration in Figure 2 specifies a list format with `<attribute>: <value(s)>` per entry. EVAPORATE outputs in this format can be de-serialize into a table.

*Prioritizing common attributes.* The list of extracted attributes and values can contain the niche attributes for specific documents, whereas a common database design principle is to capture the high frequency attributes [15]. Therefore EVAPORATE takes the union of attributes outputted across documents and ranks by frequency to enable prioritizing head attributes.

**Analysis.** We analyze this *direct extraction* implementation, EVAPORATE-DIRECT, along the axes of our three desiderata. Results processing the documents with EVAPORATE-DIRECT are reported in Table 3 and are discussed in detail in Section 4.

Overall, the quality matches or exceeds the baseline systems (described in Section 4), on 8 of the 16 settings. This is surprising given the simplicity — i.e. EVAPORATE-DIRECT uses *one* fixed prompt to process *all 16 settings*. However, fundamental cost limitations impede the real-world deployment of this approach.

However, the high cost of this implementation limits its applicability to large, recurring workloads. The number of tokens processed by the LLM scales linearly with the size of the data lake, $O(n)$. Data lakes can contain *billions* of documents [35, 50]. Further, in most organizations, data processing is not a one time cost. Data lakes are dynamically changing, so EVAPORATE-DIRECT would need to be *repeatedly* applied.

## 3.2 EVAPORATE-CODE

In this section, we present EVAPORATE-CODE, which significantly reduces cost compared to EVAPORATE-DIRECT. Here, we perform schema identification separately from value extraction, which allows us to exploit fundamental differences between the sub-tasks to reduce cost. In schema identification, we find that we only need to process a small sample of documents because attributes are consistent across documents. On the other hand, in order to extract values, we must process every document. However, the ways in which values appear across documents (*i.e.* their relative positions in the document) tend to be consistent, meaning the extraction logic is consistent across documents.

The two steps of the decomposed implementation are:

(1) **Schema synthesis.** (Section 3.2.1) We observe that the attribute outputs contain relatively consistent `<attributes>`, even though the `values` differ from document to document. To exploit this redundancy, EVAPORATE-DIRECT prompts an LLM to analyze a small sample of documents to identify attributes for the output schema For example, given a sample of the Medical Device FDA Reports, the LLM outputs a devices table with attributes like `"510(k) number"`.

(2) **Function synthesis** (Section 3.2.2). We observe consistencies in how attributes are embedded across documents. E.g., the `510(k)` code in the FDA documents always starts with the letter "k" and the `player position` attribute is always in the HTML "infobox" element in NBA player Wiki
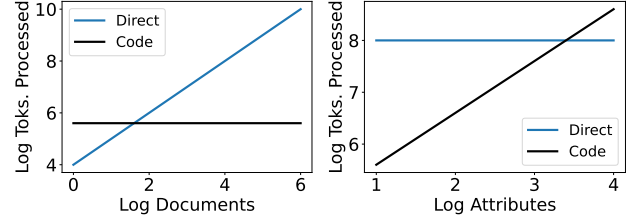


Figure 3: Tradeoffs between processing the documents via direct prompting (Direct) versus code synthesis (Code). For small data lakes and large numbers of attributes, Direct is sufficient. As the number of documents grows, Code is orders-of-magnitude more efficient. Left is evaluated at 10 attributes, Right at 10K documents, assuming 10K tokens per document.

pages. A researcher would likely exploit such redundancies when manually scraping the documents for analysis. In EVAPORATE-CODE, we propose to use the LLM to automatically synthesize a data-lake-specific suite of *functions*, that can then be applied at scale to process many documents.

Next, we provide details for each sub-task.

*3.2.1 Schema Synthesis .* EVAPORATE first uses an LLM to identify attributes $A = \{a_1, a_2, ...a_m\}$ for the output schema.

*Generating candidate attributes* Concretely, we sample a set $\tilde{D}$ of $k << n$ documents from $D$. For each, we prompt the LLM to extract the most useful attributes from the document as in EVAPORATE-DIRECT. Recall this yields a set of attributes ranked by how frequently they were extracted across documents. We retain attributes that are explicitly mentioned in the document to ensure provenance in schema identification.

*Re-ranking candidate attributes* Because EVAPORATE now identifies the attributes from a small set of documents, we observe that EVAPORATE's ranking is noisier than when every document was processed in EVAPORATE-DIRECT, i.e. an important attribute may be selected by the LLM a few times amongst the $k$ documents. Thus, we show the LLM a union of extracted attributes and prompt it to identify the most useful attributes (prompt in [5]). The frequency-based rank is upweighted if the attribute is in the LLM output.

*3.2.2 Function Synthesis.* Given the attributes $A = \{a_1, a_2...a_m\}$, the objective of EVAPORATE-CODE's second phase is to extract the values of the attributes for each document $d_i \in D$. Our key insight, as discussed, is that attribute-values are expressed in similar ways from document to document. To exploit this, instead of processing every document with the LLM to extract values for attribute $a_i$, we propose to use the LLM to *generate code* that can then be reused to process many documents.

Figure 4 shows an EVAPORATE function synthesis prompt. The in-context examples show pairs of text snippets and functions to extract an attribute of interest. EVAPORATE searches the data lake via a simple keyword search for document portions that mention $a_i$, and includes this in the prompt. EVAPORATE synthesizes functions for attributes following the rank-order of attributes derived during schema synthesis. This means that values of the most relevant (and
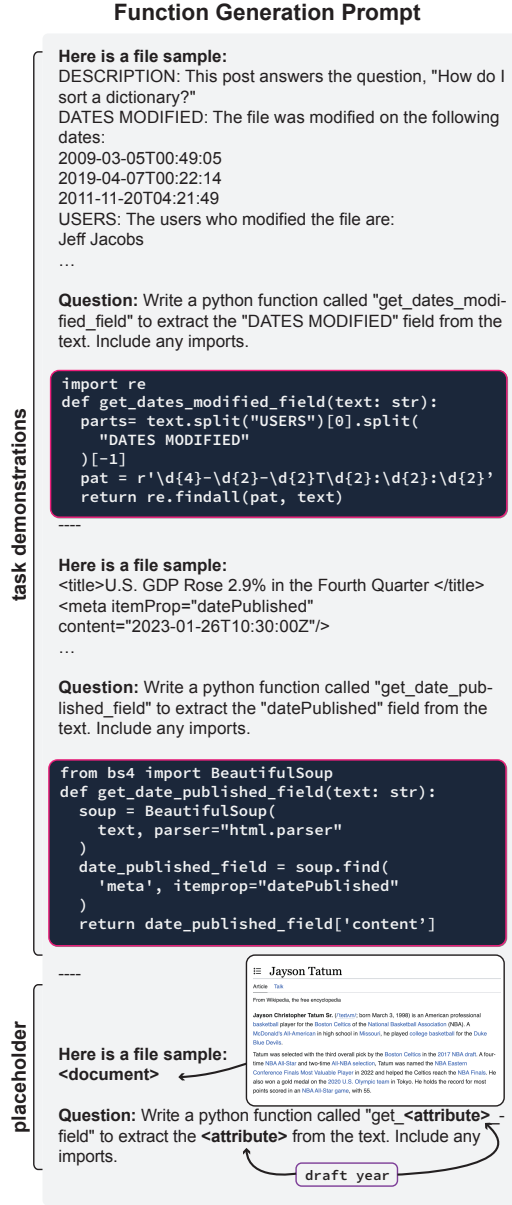
**Function Generation Prompt**

**Here is a file sample:**
DESCRIPTION: This post answers the question, "How do I sort a dictionary?"
DATES MODIFIED: The file was modified on the following dates:
2009-03-05T00:49:05
2019-04-07T00:22:14
2011-11-20T04:21:49
USERS: The users who modified the file are:
Jeff Jacobs
…

**Question:** Write a python function called "get_dates_modi-fied_field" to extract the "DATES MODIFIED" field from the text. Include any imports.

```python
import re
def get_dates_modified_field(text: str):
    parts= text.split("USERS")[0].split(
        "DATES MODIFIED"
    )[-1]
    pat = r'\d{4}-\d{2}-\d{2}T\d{2}:\d{2}:\d{2}'
    return re.findall(pat, text)
```

----

**Here is a file sample:**
<title>U.S. GDP Rose 2.9% in the Fourth Quarter </title>
<meta itemProp="datePublished"
content="2023-01-26T10:30:00Z"/>
…

**Question:** Write a python function called "get_date_pub-lished_field" to extract the "datePublished" field from the text. Include any imports.

```python
from bs4 import BeautifulSoup
def get_date_published_field(text: str):
    soup = BeautifulSoup(
        text, parser="html.parser"
    )
    date_published_field = soup.find(
        'meta', itemprop="datePublished"
    )
    return date_published_field['content']
```

----

**placeholder**

**Here is a file sample:**
<document>

**Question:** Write a python function called "get_<attribute>_-field" to extract the <attribute> from the text. Include any imports.

draft year

**Figure 4: A representative prompt for function synthesis, containing two data lake agnostic in-context examples.**

frequent) attributes as determined by EVAPORATE are extracted first. The user can stop the synthesis when desired.

**Analysis.** We briefly analyze the EVAPORATE-CODE implementation along the axes of our three desiderata. Results processing the documents with EVAPORATE-DIRECT are reported in Table 3 and are discussed in detail in Section 4.

*Cost.* Figure 3 demonstrates the asymptotic differences in cost between EVAPORATE-DIRECT and EVAPORATE-CODE. EVAPORATE-CODE is asymptotically more efficient as a function of the number

of documents: the number of LLM calls required with function generation is proportional to the number of attributes, not the number of documents. The crossover point is at ~40 documents. Meanwhile, EVAPORATE-DIRECT has the potential to extract multiple attributes from the in-context document per inference call, while EVAPORATE-CODE requires generating new functions for each attribute. Thus, the cost of EVAPORATE-CODE grows with the number of attributes, while the cost of EVAPORATE-DIRECT approach is constant. The crossover point is at ~2,500 attributes (Figure 3).

*Quality.* The tables generated by EVAPORATE-CODE are on average 21.9 pair F1 points worse than those produced using EVAPORATE-DIRECT on the SWDE datasets (Table 2). This suggests that there is a cost-quality tradeoff between the two implementations, since EVAPORATE-CODE is much cheaper.

### 3.3 EVAPORATE-CODE+

In this section we discuss an extension of EVAPORATE-CODE, which enables significant quality improvements while keeping costs low. This implementation, which we call EVAPORATE-CODE+, synthesizes *many* candidate functions and ensembles their extractions using weak supervision. We decompose the task into three parts:

(1) **Schema identification.** (Section 3.2.1) Same as in EVAPORATE-CODE.
(2) **Function synthesis.** (Section 3.3.1) Same as in EVAPORATE-CODE, except instead of generating a single function per attribute, we generate many *candidate functions*. Below we describe techniques to encourage diversity among candidates.
(3) **Function Aggregation.** (Section 3.3.2) The synthesized candidate functions have varying qualities and coverages, making them unreliable. We then introduce a weak supervision (WS) based algorithm to aggregate over their different predictions for the attribute values across documents.

*3.3.1 Synthesizing Diverse Candidate Functions.* We find that the quality of LLM-generated functions varies significantly depending on the document chunk and in-context examples used in the prompts. To address the variability in function quality, we adopt the strategy we previously proposed in Arora et al. [4]. This strategy curates multiple diverse prompt templates for the same task (i.e. multiple function generation prompts in the style of Figure 4) and prompts the LLM with each in turn to produce a diverse set of *function candidates* $F = \{f_1, f_2, ... f_k\}$.

EVAPORATE permits the use of multiple function generation prompts. We use $P_A$ and $P_B$ (included in [5]) in this work. $P_A$ has zero in-context examples and a task description that encourages the LLM to use regex. $P_B$ has two in-context examples and a task description that encourages the LLM to import and use any Python library. We find that neither consistently outperforms the other. $P_A$ produces higher quality functions on 69%, 45%, 60%, 91%, and 31% of attributes on the 8 SWDE Movie, 5 SWDE University, FDA reports, Enron, and Wikipedia player pages settings respectively. Designing a single "perfect" prompt can be challenging so EVAPORATE aggregates results from multiple prompts.

*3.3.2 Aggregating Candidate Functions.* Next, we discuss how to combine the aggregations of the candidate functions.

*Background: Methods for Unsupervised Aggregation* Because we lack ground truth labels in our setting, it is not possible to directly evaluate the quality of the candidate functions. A popular unsupervised aggregation strategy is to take the Majority Vote (MV) across function outputs [67]. Formally, MV treats the functions as independent of one another and assigns equal weight to all function outputs. However, the functions are not of equal quality — over 40% of synthesized functions result in less than 25 Text F1 in extraction quality. Therefore, EVAPORATE uses weak supervision (WS), a popular standard statistical framework for modeling the accuracies and correlations between noisy sources of information without any labeled data [27, 57]. In WS, we learn a *label model* that is parametrized by the accuracies and correlations of the candidate functions. WS is widely used in industry [57].

Unfortunately, existing WS setups make the following assumptions that do not apply in our setting. The standard setup assumes *human-designed* functions while our setting uses *machine-generated* functions that output non-standardized extracted text.

(1) **Assumption 1: Functions will abstain on examples where they do not apply [27, 57].** The attribute value returned by a function for a document could be null for two reasons: (1) the attribute does not exist in the document (e.g. a Wikipedia page may be missing a college attribute since the players did not attended college) or (2) the attribute exists but the function was not sophisticated enough to extract it (e.g. the product code attribute could start with a lowercase "k" of uppercase "K" in FDA reports, but the particular function is only designed to extract for lowercase, resulting in empty strings for uppercase documents). Note that in (1), if the function outputs a value, it has low precision and we would want to ignore the function. Note that in (2), the function has high precision and ideally our system learns to utilize such functions *selectively*. Unfortunately, it is challenging to determine whether the function outputs null for reason (1) or (2) in our setting, whereas in the traditional WS setup with human-provided functions, humans specify this logic directly (e.g. "If the email has a URL, "vote" that it contains spam, otherwise abstain" [62]).

(2) **Assumption 2: Functions are correlated with the gold label $y$ at better than random performance [27, 57].** While this is reasonable when functions are human-designed, EVAPORATE uses machine-generated functions. We find 51% of generated functions are below 50 Text F1.

(3) **Assumption 3: Weak supervision is typically applied to tasks with well defined classes in a classification setting [27, 57].** In our case, the output of the functions are extracted text, and thus there is a virtually unconstrained output space of possible extractions that vary from document to document (e.g. NBA players have varied date of birth values). The *number* of unique extractions collected by functions can also differ across documents.

We propose the following approach to be able to leverage WS. Let $D_{eval}$ be a small sample of documents from the data lake $\mathcal{D}$. We have the set of generated functions $F$ and LLM $\mathcal{F}$.

**Handling function abstentions.** To estimate the probability that an empty output from a function is an abstention, we propose

to measure the fraction $e$ of the $D_{eval}$ documents for which $\mathcal{F}$ extracts a value. Intuitively, when $e$ is high, our prior should be that the attribute appears in a large fraction of documents, so we should assume functions are *abstaining* when they output empty values. When $e$ is low, the attribute appears in few documents, so we should assume the functions are *predicting* empty values. We can use $e$ to guide both our function evaluation and downstream aggregation. Note that it is possible for $\mathcal{F}$ to abstain or hallucinate values, affecting the estimate of $e$.

**Handling functions with worse than random quality.** We propose to utilize the extractions fro $\mathcal{F}$ on a small set of documents $D_{eval}$ (e.g. we use $|D_{eval}| \leq 10$) as an estimate of the ground truth extractions for those documents. We can then estimate the quality, $\hat{a}_j$, of function $f_j$ by comparing its outputs against the outputs of $\mathcal{F}$ on document $d_i \in \mathcal{D}_{eval}$. If we are in the low $e$ regime, we should evaluate the outputs on all $d \in \mathcal{D}_{eval}$. In the high $e$ regime, we should evaluate the outputs on only the $d \in \mathcal{D}_{eval}$ for which $f_j$ extracted a value. We finally filter $f_j$ if $\hat{a}_j \leq 0.5$, where 0.5 derives from the typical WS assumptions [27, 57, 65].

Note that $\mathcal{F}$ is an LLM with its own error rate $e$, affecting the estimate $\hat{a}_j$. We theoretically study the impact of $e$ on the label model learned via WS. We provide the proof in [5].

*Proposition 1: We have m functions with empirical accuracies $\hat{a}$, evaluated against noisy labels with error rate $e$, the function accuracies estimated by the weak supervision label model are $\tilde{a}$, and the measured error is below some threshold $\epsilon$. Then, if each function labels a minimum of*

$$n \geq \frac{1}{2(\gamma - \epsilon - e)} \log(\frac{2m}{\delta})$$

*datapoints, the weak supervision label model will succeed in learning accuracies such that $||a^* - \tilde{a}||_\infty < \gamma$ with a probability $1 - \delta$.*

**Handling unconstrained output spaces.** The $k$ generated functions can produce $[0..k]$ unique prediction votes for a single unlabeled document $d_i$, and the number of unique votes can differ from document $d_i$ to $d_j$. Therefore, for each $d_i \in \mathcal{D}$, we bucket the unique votes and take the $b$ buckets representing the most frequently occurring votes. The votes for functions that outputted values outside the top-$b$ are marked as abstentions. If the number of unique votes is $< b$, placeholder values are inserted into the top-$b$. Finally, as the "classes" differ across documents, we introduce a constraint to the objective function encouraging the class-conditional accuracies to be equal.

After addressing these assumptions, we can leverage prior approaches to aggregate the noisy extractions from the function candidates into higher-quality extractions as in [4, 57]. Under WS, the output of each function is viewed as a "vote" for the true label and the objective is to construct a latent graphical model to account for the varied accuracies and correlations amongst the functions, without access to any labeled data. Our aggregation method is summarized in Algorithm 1.

**Analysis.** We briefly analyze the EVAPORATE-CODE+ implementation along the axes of our three desiderata. Results processing the documents with EVAPORATE-CODE+ are reported in Table 3 and are discussed in detail in Section 4.

*Cost.* As with EVAPORATE-CODE, the number of tokens processed by the LLM in EVAPORATE-CODE+ is fixed with respect to the number of documents. Figure 3 demonstrates the asymptotic differences in cost between EVAPORATE-DIRECT and EVAPORATE-CODE. The number of tokens that must be processed by the LLM grows only by a constant factor: the number of function candidates generated. The user can set this number to balance cost and quality.

*Quality.* Of the three implementations, EVAPORATE-CODE+ produces the highest quality tables. EVAPORATE-CODE+ outperforms EVAPORATE-DIRECT by 12.1 F1 points (22%) on average, while using far fewer computational resources. Using function aggregation leads to an improvement of 25.1 F1 points over EVAPORATE-CODE.

## 4 EVALUATIONS

We now evaluate EVAPORATE, validating the following claims:

- **Function synthesis enables asymptotic cost reductions for processing data with LLMs.** There has been significant recent interest in developing various data management applications with LLMs [17, 36, 40, 49]. Prior work directly processes data with the LLM. EVAPORATE-CODE+ reduces the number of tokens the LLM needs to process by 110x relative to EVAPORATE-DIRECT.
- **Function synthesis + aggregation results in higher quality than direct extraction.** Despite the fact that EVAPORATE-DIRECT processes each document with the LLM directly, EVAPORATE-CODE+ performs 10.1 F1 points (18%) better on average. Based on comparisons with EVAPORATE-CODE, which only synthesizes one function, we show that function aggregation is key in enabling the improvements.
- **EVAPORATE achieves higher quality than state-of-the-art baselines, while exposing a more general interface.**

EVAPORATE-CODE+ expresses tasks via merely *six* natural language prompts (all provided in [5]) and uses no training. Yet, it exceeds SoTA systems by 3.2 F1 (6%) points when generating tables from scratch and 6.7 points (10%) when extracting pre-defined gold attributes. Meanwhile, it supports a broader range of settings than any of these baselines.

- **The identified tradeoffs hold across language models.** We evaluate on four models from three unique providers [6, 43, 52]. We find EVAPORATE-DIRECT and EVAPORATE-CODE+ remain competitive in quality across LLMs.

### 4.1 Experimental Setup

We primarily evaluate EVAPORATE on the end-to-end task of *structured view generation*. For the purpose of comparison to prior work, we also evaluate on the sub-task of *closed information extraction*. We first define these tasks, their metrics, and the baselines. We then provide implementation details for EVAPORATE.

**Structured view generation task.** This captures the end-to-end task of identifying the schema and populating the output table. This task is often discussed as a vision system [16], and given the difficulty of this task, there are limited comparable works. We therefore compare to the closest line of work, OpenIE systems, where the task is to extract all facts from documents [7, 51]. We compare to two sets of baselines: (1) Deng et al. [22], Lockard et al. [45, 46] for HTML-specific OpenIE, and (2) Kolluru et al. [39] for generic unstructured text. The former models explicitly use the HTML-DOM tree structure to process the page, assuming attribute values are leaf nodes, and explicitly train on documents from the domain of interest. The latter class of systems first label sentences using linguistic tools (*i.e.* dependency parsers, part of speech taggers, and named entity taggers), and fine tune LLMs over these features to perform the task [71].

*Metrics.* The standard metric is Pair F1 [22, 45], an F1 score applied to the predicted vs. gold sets of tuples of the form (document ID $d_i$, attribute $a_j$, value $r_{i, j}$). The tuple must exactly match a tuple in the ground truth to be marked correct. Since EVAPORATE ranks the attributes and generates functions in this order, for fair comparison, we report OpenIE scores for all tuples up to $k$ attributes, where $k$ is the number of gold attributes for the setting. We note that the prior systems extract all-or-no tuples, in contrast.

**Closed information extraction task.** This captures the setting where the user provides a pre-defined schema and EVAPORATE is used to populate the table. We compare to state-of-the-art approaches for closed IE including: (1) Deng et al. [22], Lockard et al. [45, 46] for HTML-specific ClosedIE and (2) Clark et al. [19], He et al. [33] for generic unstructured text. The former models explicitly use the HTML-DOM tree structure to process the page, assuming attribute values are leaf nodes, and explicitly train on documents from the test domain. The latter are pretrained LLMs that have been fine tuned on massive amounts of labeled (`attribute, value`) pairs [56]. We report ClosedIE results using the Text F1 metric on a value-by-value basis across each document.

**EVAPORATE Implementation Details.** In the following experiments, we instantiate EVAPORATE with currently popular, LLM APIs. Experiments in Sections 4.3 and 4.4.1 use `text-davinci-003` from OpenAI. In Section 4.4.2, we evaluate additional LLMs from three

model providers. For experiments, we use 10 sample documents per data lake for the schema synthesis, function synthesis, and function verification. We apply Algorithm 1 over the top-10 scoring functions that are synthesized for each attribute and data lake. The prompts remain constant across data lakes and models. In [5], we provide ablations that show how the system's quality changes as we vary the number of sample documents and top-$k$ functions.

When the measuring cost for alternate implementations of Evaporate, we compute total number of tokens processed by the LLM to perform the end-to-end task (*i.e.* the sum of the number of tokens in the prompt and model generation). We use this metric because the wall-clock time and dollar cost of a model fluctuate, but both should be proportional to the number of tokens processed.

## 4.2 Evaluation Settings

We evaluate Evaporate on 16 settings representing a range of real-world data lakes. First, we use a benchmark suite of 13 Movie and University websites to compare Evaporate to state-of-the-art information extraction systems [22, 32, 45]. Next, to evaluate on more unstructured data (*i.e.* non-HTML) we turn to: **Enron** a corporate email corpus that has been analyzed in over three thousand academic papers [3, 34, 37], **FDA 510(k)** reviews for premarket notification submissions for medical devices, which have been the subject of multiple important research studies [68, 72], and **NBA** Wikipedia pages for NBA players, which include more complex HTML than the existing benchmarks [22]. We release the benchmarks and provide additional details in [5]. Here we briefly describe the properties we aim to study with each setting:

(1) **Benchmark Suite: SWDE Movies & Universities** SWDE is the standard benchmark for document-level IE in prior work [22, 32, 45, 46, inter alia.]. There are 8 sets of webpages for Movies (e.g. IMDB) and 5 sets of webpages for Universities (e.g. US News). For each website, the benchmark contains 1063-2000 pages and annotations for 8-274 attributes. We use SWDE to compare to the state-of-the-art and test on a range of attribute types, e.g. simpler Movie "runtime" through complex Movie "cast" and popular Movie "director" through infrequent "second assistant director".

(2) **Complex HTML: NBA** As SWDE attributes always occur in separate leaf nodes of the HTML-DOM tree, we use NBA Player Wikipedia pages to evaluate on more complex HTML. E.g., the NBA draft attribute contains the draft round, year, pick number, and team by which the player was selected. We evaluate on 100 randomly selected player pages (spanning the 1940s-present) and 19 attribute annotations.

(3) **Unstructured Text: Enron and FDA** We observe a lack of existing benchmarks for document-level IE over unstructured text — intuitively, this setting has been challenging with prior generations of models due to the lack of *any* grounding structure whatsoever (*i.e.* recall current systems rely on HTML-DOM elements or sentence-level NER, dependency, and POS tags). We turn to the Enron and FDA settings described above. The Enron setting contains 15 gold attributes and 500k documents. The FDA setting contains 16 gold attributes and 100 PDF documents, which are up to 20 pages long, randomly sampled from FDA 510(k).



**Figure 5: T-SNE Visualization of documents in the SWDE dataset. T-SNE is performed on the first 16 principal components of TF-IDF vectors. Colors indicate the source website.**

**Dataset Protocols** Because the baselines we compare against on SWDE require training data, they perform website-wise cross validation (*i.e.* train on some websites and evaluate on others). They do this for several combinations such that every website appears in the evaluation set. Evaporate, in contrast, does not require any training data. We simply evaluate Evaporate on all websites so that our method is evaluated on the same examples as the baselines.

We process each dataset with Evaporate separately, to match the protocol used by the baselines [22]. However, we may not have access to the sources of the documents in a real-world data lake – they may be mixed together. Using a standard TF-IDF vectorizer and K-means clustering to the mixture of documents, we verify we can perfectly recover the document sources without any labeled data or supervision (Figure 5, details in [5]). Intuitively, clustering semi-structured data may be simple due to the rich formatting.

## 4.3 Comparing Evaporate to Baselines

First we validate that Evaporate outperforms baselines, defined in Section 4.1, both in generality (i.e. the flexibility to support data from different domains and formats) and quality. We then compare the efficiency of Evaporate-Code+ vs. the baselines.

### 4.3.1 Quality and generality comparisons.

*Systems for semi-structured text.* Shown in Table 2, Evaporate outperforms the state-of-the-art on SWDE. We compare to the metrics reported the baseline works. Recall Evaporate uses no training whatsoever and can be applied across document formats (HTML, PDF, TXT). In contrast, the baselines are limited to HTML and explicitly perform supervised learning using labels from webpages within the Movie and University domains respectively [22, 46]. E.g., Deng et al. [22] assumes attribute values are the leaf-nodes of the HTML-DOM tree and thus does not work on non-HTML.

The baseline systems restrict scope to attributes that are specifically mentioned in the HTML <body> text, even though attributes are frequently mentioned in the HTML header (e.g. within <title> elements) and tags (e.g. <a href='year/2012'>). We validate Evaporate can identify and extract attributes mentioned anywhere in the document. We extend the SWDE benchmark to include the attributes scattered throughout the full HTML and find Evaporate achieves 52.2 and 49.0 on Movies and University respectively on the more challenging setting. We release the new annotations.

**Table 1: Quality of EVAPORATE-CODE+ evaluated on ClosedIE in Text F1 and OpenIE in Pair F1 using text-davinci-003.**

| Source (Format) | CLOSEDIE | OPENIE | | |
| | F1 | R | P | F1 |
| --- | --- | --- | --- | --- |
| FDA (TXT) | 80.1 | 62.0 | 68.1 | 64.9 |
| Enron Emails (TXT) | 93.3 | 80.3 | 94.6 | 86.9 |
| Wiki NBA (HTML) | 84.7 | 55.7 | 88.2 | 68.2 |
| SWDE Movie (HTML) | 79.5 | 48.5 | 71.0 | 56.8 |
| SWDE University (HTML) | 73.7 | 50.9 | 71.4 | 59.0 |
| **Average** | **82.3** | **58.9** | **78.5** | **66.7** |

**Table 2: Comparisons to state-of-the-art on ClosedIE in Text-F1 and OpenIE in Pair F1. The baselines train on in-domain documents, while EVAPORATE uses no training [22].**

| System | SWDE MOVIE | | SWDE University | |
| | Closed | Open | Closed | Open |
| --- | --- | --- | --- | --- |
| ZeroShot Ceres [46] | - | 50.0 | - | 50.0 |
| RoBERTa-Base | 49.3 | 35.6 | 36.6 | 38.0 |
| RoBERTa-Structural | 47.7 | 39.9 | 46.5 | 42.3 |
| DOM-LM [22] | 71.9 | 54.1 | 68.0 | 55.2 |
| EVAPORATE-DIRECT | **84.4** | 45.2 | 72.6 | 53.8 |
| EVAPORATE-CODE | 55.0 | 33.0 | 40.5 | 22.2 |
| EVAPORATE-CODE+ | 79.5 | **56.8** | **73.7** | **59.0** |

*Systems for unstructured text.* We are not aware of strong baselines that apply beyond HTML document formats. The most relevant baseline is the OpenIE6 system for performing OpenIE over any unstructured text from Kolluru et al. [39]. We find the system only handles well formatted sentences and struggles to extend to heterogeneous data types. We find that even when documents contain full sentences, the system extracts an extremely large set of relations and does not enforce consistent extractions across documents. For instance, on a sample FDA 510(k) document, OpenIE6 extracts 427 relations with 184 relations having a confidence level at 0.99. We include a detailed error analysis in [5].

*4.3.2 Efficiency comparisons.* We compare the efficiency of EVAPORATE with the (estimated [28]) 175B parameter OpenAI model vs. the baseline, which uses a pretrained 125M parameter RoBERTa model [22], decomposed in terms of pretraining, fine-tuning, inference, and parameter (memory) cost. Using FLOPS as reported in the OpenAI reference work Brown et al. [14], for data settings with $n$ documents and $m$ attributes, the costs are:

- **RoBERTa** The model is 125M parameters, total pretraining FLOPS is 1.50E+21, and inference FLOPS per token is 2 × Number of Parameters, which is 0.250 GFLOPS. The total inference cost is computed by

$$n \times \frac{tokens}{document} \times 0.250 \text{ GFLOP}$$

- **GPT3-175B** The model is 175B parameters, total pretraining FLOPS is 3.14E+23, and inference FLOPS per token is 2 × Number of Parameters, which is 350 GFLOPS. The total inference cost is computed by

$$m \times P \times \frac{tokens}{chunk} \times 350 \text{ GFLOP}$$

where $P$ is the number of prompts per attribute. Note that for our evaluated implementation $P \approx 10c$ since function generation is performed on 10 documents.

EVAPORATE uses a model with 1,400x more parameters and 300x higher pretraining cost. However, users of the baselines likely need to locally fine-tune and host the models. The inference costs of the baseline and EVAPORATE on our datasets are in the same order of magnitude (summarized in [5]). The extended cost comparison and analysis are provided in [5]. Users should select a method depending on the data setting, i.e. the number of documents and attributes. We note that EVAPORATE allows users to tradeoff quality and efficiency by changing the underlying language model to smaller variants.

## 4.4 Comparing Implementations of EVAPORATE

This work proposes a fundamental tradeoff space between directly processing data workloads with LLMs vs synthesizing code that does the processing. We first discuss the tradeoffs for a fixed LLM (text-davinci-003), which is the current best-in-class LLM [42] (Section 4.4.1), and next across a range of LLMs trained by three distinct model providers (Section 4.4.2).

*4.4.1 Tradeoffs between EVAPORATE Implementations.* As detailed in Section 3.2, the base routine ("EVAPORATE-DIRECT") in EVAPORATE entails directly processing documents with the LLM, while the optimized routine ("EVAPORATE-CODE") synthesizes functions for processing. Next we evaluate these along our desiderata.

**Generality is maintained.** LLMs take text as input and provide text as output — this unified natural language interface means EVAPORATE-DIRECT and EVAPORATE-CODE can ingest any document format without additional engineering. *Critically, our results with EVAPORATE require no user effort, no training whatsoever, and no customization when applied to the 16 different settings.*

**Asymptotic cost reduction.** Figure 3 demonstrates the asymptotic differences in cost between directly processing the data lake with EVAPORATE-DIRECT vs. with EVAPORATE-CODE+. (Figure 3 Left) EVAPORATE-DIRECT is asymptotically more efficient as a function of the number of documents in the data lake. The number of LLM calls required with function generation is proportional to the number of attributes to be extracted, not the number of documents. The crossover point is at ~40 documents.

(Figure 3 Right) EVAPORATE-DIRECT can extract multiple (*i.e.* every) attribute in the in-context documents in a single inference call, while EVAPORATE-CODE+ synthesizes new functions for each attribute. Thus, the cost of function synthesis grows with the number of attributes, while the cost of EVAPORATE-DIRECT is constant. The crossover point is at ~2,500 attributes.

Empirically across our settings, EVAPORATE-CODE+ realizes a 110x average reduction in the number of tokens the LLM needs to process (assuming 10k documents per setting and 378× given the true benchmark sizes) in the number of tokens the LLM must process compared to EVAPORATE-DIRECT (Table 3). Further, data lakes are constantly changing and functions can be reused while EVAPORATE-DIRECT would need to be re-run, multiplying the cost.

In runtime, we observe that the generated functions are efficient in processing the documents. For example, over the 9,500 function

**Table 3: Quality (OpenIE Pair F1) and cost (number of tokens processed by the LLM) for producing the structured views. We compare the direct prompting and code synthesis implementations using text-davinci-003. EVAPORATE-CODE+ is evaluated on the full datasets, while EVAPORATE-DIRECT is evaluated on a randomly sampled proportion of documents due to the cost (20% of FDA and Wiki, 2% of SWDE, and 0.0002% of Enron, given the respective sizes).**

| Source (Format) | EVAPORATE-DIRECT | | | EVAPORATE-CODE+ | | | Relative Performance | |
| | Quality F1 | Cost / 10K Documents | | Quality F1 | Cost / 10K Documents | | Quality | Cost Reduction |
| | | Tokens (M) | Cost ($) | | Tokens (M) | Cost ($) | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| FDA (TXT) | 45.5 | 145.6 | 2,900 | 62.8 | 1.9 | 38 | +17.3 | 77x |
| Enron Emails (TXT) | 93.8 | 21.2 | 425 | 86.9 | 0.6 | 12 | -6.9 | 35x |
| Wiki NBA (HTML) | 44.8 | 650.1 | 13,000 | 68.2 | 3.0 | 60 | +23.4 | 217x |
| SWDE Movie (HTML) | 45.2 | 282.9 | 5,660 | 56.8 | 2.3 | 46 | +11.6 | 123x |
| SWDE University (HTML) | 53.8 | 190.1 | 3,800 | 59.0 | 1.9 | 38 | +5.2 | 100x |
| **Average** | **56.6** | **258** | **5,157** | **66.7** | **1.9** | **39** | **+10.1** | **110x** |

runs (from 95 functions evaluated on 100 documents each) in the FDA 510(k) setting, we find that the average time to run one function over one document is 0.00025s on a 2 CPU machine.

**Improved quality and reliability.** Even though EVAPORATE-DIRECT directly processes each document with the LLM, EVAPORATE-CODE+ surprisingly performs 10.1 F1 (18%) better (Table 3).

*What are the failure modes of EVAPORATE-DIRECT?* The method yields inconsistent generations. On the Medical FDA report setting:(1) The LLM misses an average of 4.4 attributes that are present in the gold schema (27.5% of gold attributes) per document. Among the gold attributes that are missed, all *are extracted* in at least one document. (2) Further, the LLM outputs an average of 9.7 attributes or values that are not explicitly mentioned in the documents. (3) Finally, attributes are reworded in diverse ways across documents — the attribute classification is extracted in 4 different ways across the sample of 10 documents (i.e. "classification", "device classification", "regulatory information", missing). Since the error modes are quite varied, it is unclear how to improve quality.

*Why does EVAPORATE-CODE+ improve quality?* We validate that our Algorithm 1 for selecting and aggregating functions leads to the quality improvements over EVAPORATE-DIRECT.

**Synthesizing diverse functions** We find that using diverse prompts helps address the lack of reliability in function synthesis. To synthesize functions, EVAPORATE-CODE+ uses a prompt template that includes one-to-two in-context examples and a placeholder for the inference example, i.e. document text (Figure 4). We can produce multiple prompts in the template by swapping the in-context examples or sampling more documents (change the inference example). We find both means of increasing diversity benefit quality:

- **In-context demonstrations** Our implementation (Table 1) instantiates two prompts by swapping in-context demonstrations, $P_A$ and $P_B$. Quality using $P_A$ or $P_B$ alone is 8.5 and 8.0 F1 points worse than using both to synthesize functions on SWDE Movie and SWDE University respectively.
- **Inference documents** Using three versus five sample documents in the prompts for EVAPORATE-CODE+, the ClosedIE and OpenIE quality improve by 6.8 F1 points (9%) and 6.5 F1 points (14%) respectively, averaged across the 16 settings.

**Estimating function quality using the LLM.** In Table 4, we first evaluate the two unsupervised aggregation baselines in prior work off-the-shelf: Majority Vote (MV) and Weak Supervision (WS)

**Table 4: Quality under alternate approaches of aggregating the synthesized functions. Baselines are in the left columns: Majority Vote (MV) and Weak Supervision (WS). Components of Algorithm 1 are in the right columns: "Abstain" accounts for abstensions and "Filter" filters low quality functions.**

| Source | MV | WS | WS Filter | WS Abstain + Filter |
| --- | --- | --- | --- | --- |
| FDA (TXT) | 52.9 | 51.1 | 55.0 | 62.8 |
| Enron Emails (TXT) | 81.4 | 82.7 | 86.9 | 86.9 |
| Wiki NBA (HTML) | 59.5 | 64.9 | 68.4 | 68.2 |
| SWDE Movie (HTML) | 44.3 | 46.3 | 56.6 | 56.8 |
| SWDE University (HTML) | 42.7 | 43.5 | 57.3 | 59.0 |
| **Average** | **56.2** | **57.7** | **64.8** | **66.7** |

[4, 57, 67]. Next we measure the effect of filtering functions and handling abstentions as proposed in Algorithm 1.

In Table 4, we observe WS with filtering provides a consistent boost across settings compared to WS — 7.1 F1 point higher average quality and up to 13.8 F1 points on the SWDE University setting. Additionally handling abstentions leads to a 1.9 F1 point increase in average quality over WS with filtering, with up to 7.8 F1 points on the FDA setting. Qualitatively, accounting for abstentions is helpful when attributes are expressed in diverse ways across documents, which is not applicable to all settings such as Enron. These results highlight the importance of EVAPORATE-CODE+'s aggregation approach for the system's overall reliability. Without Algorithm 1, quality does not improve over EVAPORATE-DIRECT.

*4.4.2 Understanding the Tradeoff Space across Varied Language Models.* The are an increasing number of LLMs being available. These models are trained by various providers each using distinct protocols [42]. To understand whether the tradeoffs we identified hold for different LLMs, we evaluate EVAPORATE using three additional LLMs from three different providers: (1) **GPT-4** [52], (2) **Anthropic Claude-V1** [6], and (3) **Jurassic Jumbo-2-Instruct** [43]. Results are summarized in Table 5.

*Overall results.* The quality with gpt-4 is comparable to that obtained using text-davinci-003. Both the EVAPORATE-DIRECT and EVAPORATE-CODE+ quality decrease with claude and jumbo, consistent with the results of large-scale benchmarking efforts [42], however the *relative* quality of the two implementations are similar to Table 3. Both appear to remain competitive in quality and the quality of the approaches appear to increase together.

**Table 5: OpenIE (Pair F1) results evaluating EVAPORATE using alternate LMs from three model providers. For cost reasons, we apply EVAPORATE-DIRECT to samples of 10 documents each. For fair comparison, we report the score of EVAPORATE-CODE+ on the same sample instead of the full set of documents. $k$ is the number of gold attributes for the setting.**

| Source (Format) | EVAPORATE-DIRECT | | | | | EVAPORATE-CODE+ | | | | | SCHEMA ID |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | FDA | Wiki | Movie | University | Enron | FDA | Wiki | Movie | University | Enron | F1@k |
| OpenAI GPT-4 [52] | 59.2 | 40.5 | 35.1 | 56.1 | 92.7 | 57.5 | 61.4 | 54.9 | 57.2 | 85.5 | 67.3 |
| Anthropic Claude-V1 [6] | 45.1 | 20.6 | 27.5 | 44.3 | 88.1 | 44.4 | 33.5 | 38.7 | 30.4 | 84.7 | 69.0 |
| Jurassic Jumbo-2-Instruct [43] | 25.9 | 0.0 | 13.3 | 29.2 | 90.3 | 1.2 | 0.0 | 20.6 | 18.6 | 85.7 | 62.3 |

We find the precision of EVAPORATE-CODE+ remains high across models. Algorithm 1 helps EVAPORATE filter the low quality functions and if this eliminates all the candidate functions, the attribute is excluded from the output table. We find that when an attribute *is* included in the output, it has high precision, consistent with Table 1 where the precision with `text-davinci-003` is almost 20 points higher than the recall. The average precision scores corresponding to EVAPORATE-CODE+ in Table 5 are 70.9 (gpt-4), 67.6 (claude), and 50.9 (jumbo) using EVAPORATE-CODE+ and in contrast are 61.9 ( gpt-4), 55.1 (claude), and 49.9 (jumbo) using EVAPORATE-DIRECT, emphasizing a precision-recall tradeoff between approaches.

*Understanding the errors.* Overall, EVAPORATE relies on versatile reasoning capabilities (*i.e.* to identify the schema and extract attribute values *directly* from *noisy* provided context, and the ability to synthesize code) and excitingly, the results validate that these capabilities co-exist within multiple model families. We investigate which of the required reasoning capabilities contributes to lower quality in comparison to `text-davinci-003`. We find that the schema synthesis step plays a small role. Considering the top ranked schema attributes according to EVAPORATE, we measure the average F1@k between the predicted and gold sets of attributes, where $k$ is the number of gold attributes per setting. The average F1@k for `text-davinci-003` is 71.9, and the right-hand column of Table 5 shows the alternate models perform comparably.

We find the two main sources of errors are (1) the inability to generate a function for particular attributes, and (2) occasionally, low quality *direct* extractions in particular cases (e.g., `claude` may respond "I'm not sure, please give me more information." in a Chat-Bot style, when prompted to extract an attribute value). The models we evaluate are optimized for ChatBot applications [38].

## 5 RELATED WORK

*Structured Querying of Heterogeneous Data.* Converting heterogeneous data to structured databases is a long standing data management problem [12, 16, 31, inter alia.]. In contrast to systems for knowledge base construction (KBC) or closed information extraction (IE) [61, 70], which assume there is a predefined schema and focus on populating the database according to the schema, the setup we focus on relies on OpenIE. OpenIE is the task of extracting useful facts without access to a predefined ontology (i.e. the types or categories of facts to be extracted) [7, 20]. Given the breadth of input documents, the ability to construct a schema and populate the corresponding database on-the-fly is useful.

Existing systems for this problem introduce assumptions about the data-domain [22, 55, 64], file-format (e.g., XML files) [30], or the syntactic patterns of useful facts [12, 15, 25, 31, 39, 48, 51, 71]. For instance, in early systems, Cafarella et al. [15] focuses on facts

expressed as triples (two entities with a descriptive string of their relationship in between) with hypernym "is-a" relationships between the entities. The recent deep learning based systems (1) require domain- and document-format specific training, (2) focus on reasoning over sentences, in contrast to long documents, and (3) rely on high quality linguistic tools (e.g. dependency parse, POS, NER) to help introduce structure over unstructured text [39, 71].

For the narrower problem of generating structured views from *web data* [15, 23, 25], the current state-of-the-art approaches use (1) distant supervision to train site-specific extraction models [45] (**domain specific**), and (2) rely on assumptions about where in the HTML-DOM attributes and values are located (**format specific**) Deng et al. [22], Lockard et al. [46]. We investigate the feasibility of a domain and document-format agnostic approach.

*Language Models for Data Management.* Given the recency of in-context learning, there are few works exploring the benefits for data processing. Most closely related, Chen et al. [17] presents a system for querying heterogeneous data lakes with in-context learning. The proposed approach involves processing *every document* with the LLM to extract values of interest. We propose an alternate approach and tradeoff space for processing data with LLMs.

Other recent work applies language models for tasks such as data wrangling [49] or code generation for SQL queries [63]. Unlike EVAPORATE, supporting various data formats, these prior approaches focus on relational data only and design manual prompts to demonstrate high quality.

*Data Programming.* We build on work in data programming and weak supervision [57]. EVAPORATE automatically generates functions rather than using human-designed functions. We use WS on open ended tasks in contrast to the classification tasks considered in the prior work on automated WS [9, 65]. We show how to use the LLM to handle abstentions and filter low quality functions.

## 6 CONCLUSION

We propose EVAPORATE, a system that uses LLM in-context learning to generate structured views of semi-structured data lakes. We identify and explore a cost-quality tradeoff between processing data directly with an LLM versus synthesizing and aggregating multiple code snippets for data processing. We present an algorithm and theoretical analysis for applying weak supervision to aggregate the code snippets. The code-based approach aims to exploit the structural redundancies that occur in corpora of semi-structured documents. We validate EVAPORATE on 16 unique data settings spanning 5 domains and 3 document formats, considering the cost, quality, and generality of the system. Our study highlights the promise of LLM-based data management systems.

# REFERENCES

[1] April 2023. Wikipedia Statistics. https://en.wikipedia.org/wiki/Special:Statistics

[2] Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. 2022. Large Language Models are Few-Shot Clinical Information Extractors. *The 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (2022).

[3] Simran Arora, Patrick Lewis, Angela Fan, Jacob Kahn, and Christopher Ré. 2023. Reasoning over Public and Private Data in Retrieval-Based Systems. *Transactions of Computational Linguistics (TACL)* (2023).

[4] Simran Arora, Avanika Narayan, Mayee F. Chen, Laurel Orr, Neel Guha, Kush Bhatia, Ines Chami, Frederic Sala, and Christopher Ré. 2023. Ask Me Anything: A simple strategy for prompting language models. *International Conference on Learning Representations (ICLR)* (2023).

[5] Simran Arora, Brandon Yang, Sabri Eyuboglu, Avanika Narayan, Andrew Hojel, Immanuel Trummer, and Christopher Ré. 2023. Language Models Enable Simple Systems for Generating Structured Views of Heterogeneous Data Lakes. (2023). https://www.dropbox.com/scl/fi/3gt3ixdbvp986ptyz5j4t/VLDB_Revision.pdf?rlkey=mxi2kqp7rqx0frm9s7bpttwcq&dl=0

[6] Amanda Askell, Yushi Bai, Anna Chen, Dawn Drain, Deep Ganguli, T. J. Henighan, Andy Jones, and Nicholas Joseph et al. 2021. A General Language Assistant as a Laboratory for Alignment. *arXiv:2112.00861v3* (2021).

[7] Michele Banko, Michael J. Cafarella, Stephen Soderland, Matthew G Broadhead, and Oren Etzioni. 2007. Open information extraction from the web. *IJCAI* (2007).

[8] David W Bates, David M Levine, Hojjat Salmasian, Ania Syrowatka, David M Shahian, Stuart Lipsitz, Jonathan P Zebrowski, Laura C Myers, Merranda S Logan, Christopher G Roy, et al. 2023. The Safety of Inpatient Health Care. *New England Journal of Medicine* 388, 2 (2023), 142–153.

[9] Benedikt Boecking, Willie Neiswanger, Eric Xing, and Artur Dubrawski. 2021. Interactive weak supervision: Learning useful heuristics for data labeling.

[10] Rideout J. R. Dillon M. R. Bokulich N. A. Abnet C. C. Al-Ghalith G. A. Alexander H. Alm E. J. Arumugam M. et al. Bolyen, E. 2019. Reproducible, interactive, scalable and extensible microbiome data science using qiime 2. In *Nature biotechnology.*

[11] Rishi Bommasani, Drew A. Hudson, E. Adeli, Russ Altman, Simran Arora, S. von Arx, Michael S. Bernstein, Jeanette Bohg, A. Bosselut, Emma Brunskill, and et al. 2021. On the opportunities and risks of foundation models. *arXiv:2108.07258* (2021).

[12] S. Brin. 1998. Extracting patterns and relations from the WorldWide Web. In *WebDB.*

[13] Mirko Bronzi, Valter Crescenzi, Paolo Merialdo, and Paolo Papotti. 2013. Extraction and integration of partially overlapping web sources. *PVLDB* (2013).

[14] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

[15] Michael J. Cafarella, Christopher Re, Dan Suciu, Oren Etzioni, and Michele Banko. 2007. Structured Querying of Web Text. In *Conference on Innovative Data Systems Research (CIDR).*

[16] Michael J Cafarella, Dan Suciu, and Oren Etzioni. 2007. Navigating Extracted Data with Schema Discovery.. In *WebDB.* 1–6.

[17] Zui Chen, Zihui Gu, Lei Cao, Ju Fan, Sam Madden, and Nan Tang. 2023. Symphony: Towards Natural Language Query Answering over Multi-modal Data Lakes. *CIDR* (2023).

[18] Eric Chu, Akanksha Baid, Ting Chen, AnHai Doan, and Jeffrey Naughton. 2007. A Relational Approach to Incrementally Extracting and Querying Structure in Unstructured Data. In *VLDB.*

[19] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA: pre-training text encoders as discriminators rather than generators. In *International Conference on Learning Representations (ICLR).*

[20] W. Cohen. 2004. Information extraction and integration: An overview. *IJCAI* (2004).

[21] Lei Cui, Furu Wei, and Ming Zhou. 2022. Neural Open Information Extraction. (2022).

[22] Xiang Deng, Prashant Shiralkar, Colin Lockard, Binxuan Huang, and Huan Sun. 2022. DOM-LM: Learning Generalizable Representations for HTML Documents. (2022).

[23] Oren Etzioni, Michele Banko, Stephen Soderland, and Daniel S Weld. 2008. Open information extraction from the web. *Commun. ACM* 51, 12 (2008), 68–74.

[24] Oren Etzioni, Michael Cafarella, Doug Downey, Stanley Kok, Ana-Maria Popescu, Tal Shaked, Stephen Soderland, Daniel S Weld, and Alexander Yates. 2004. Web-scale information extraction in knowitall: (preliminary results). In *Proceedings of the 13th international conference on World Wide Web.* 100–110.

[25] Oren Etzioni, Michael Cafarella, Doug Downey, Ana-Maria Popescu, Tal Shaked, Stephen Soderland, Daniel S. Weld, and Alexander Yates. 2004. Unsupervised named-entity extraction from the Web: An experimental study. In *AAAI.*

[26] J. H. Faghmous and V Kumar. 2014. A big data guide to understanding climate change: The case for theory-guided data science. In *Big data.*

[27] Daniel Fu, Mayee Chen, Frederic Sala, Sarah Hooper, Kayvon Fatahalian, and Christopher Re. 2020. Fast and Three-rious: Speeding Up Weak Supervision with Triplet Methods. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Vol. 119. PMLR, 3280–3291.

[28] Leo Gao. 2021. On the Sizes of OpenAI API Models. https://blog.eleuther.ai/gpt3-model-sizes/

[29] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2021. The Pile: An 800GB Dataset of Diverse Text for Language Modeling.

[30] Minos Garofalakis, Aristides Gionis, Rajeev Rastogi, Sridhar Seshadri, and Kyuseok Shim. 2000. XTRACT: A system for extracting document type descriptors from XML documents. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data.* 165–176.

[31] Eugene Agichtein Luis Gravano. 2000. Snowball: Extracting Relations from Large Plain-Text Collections. In *DL '00: Proceedings of the fifth ACM conference on Digital libraries.*

[32] Qiang Hao, Rui Cai, Yanwei Pang, and Lei Zhang. 2011. From one tree to a forest: a unified solution for structured web data extraction. *SIGIR* (2011).

[33] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. DeBERTa: Decoding-Enhanced BERT With Disentangled Attention. In *International Conference on Learning Representations.*

[34] Nathan Heller. 2017. What the Enron E-mails Say About Us. https://www.newyorker.com/magazine/2017/07/24/what-the-enron-e-mails-say-about-us

[35] Nick Huss. 2023. How Many Websites Are There in the World?

[36] Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. Demonstrate-Search-Predict: Composing Retrieval and Language Models for Knowledge-Intensive NLP. *arXiv preprint arXiv:2212.14024* (2022).

[37] B. Klimt and Y. Yang. 2004. Introducing the enron corpus. In *Proceedings of the 1st Conference on Email and Anti-Spam (CEAS).*

[38] Jan Kocoń, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielaniewicz, Marcin Gruza, Arkadiusz Janz, Kamil Kanclerz, et al. 2023. ChatGPT: Jack of all trades, master of none. *arXiv preprint arXiv:2302.10724* (2023).

[39] Keshav Kolluru, Vaibhav Adlakha, Samarth Aggarwal, Mausam, and Soumen Chakrabarti. 2020. OpenIE6: Iterative Grid Labeling and Coordination Analysis for Open Information Extraction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP).*

[40] Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott Wen tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2022. DS-1000: A Natural and Reliable Benchmark for Data Science Code Generation. *ArXiv* abs/2211.11501 (2022).

[41] Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022. A Survey on Retrieval-Augmented Text Generation. *CoRR* abs/2202.01110 (2022). arXiv:2202.01110 https://arxiv.org/abs/2202.01110

[42] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, and more. 2022. Holistic Evaluation of Language Models. *ArXiv* abs/2211.09110 (2022).

[43] Opher Lieber, Or Sharir, Barak Lenz, and Yoav Shoham. 2021. Jurassic-1: Technical details and evaluation. (2021).

[44] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR* abs/1907.11692 (2019). arXiv:1907.11692 http://arxiv.org/abs/1907.11692

[45] Colin Lockard, Prashant Shiralkar, and Xin Luna Dong. 2019. OpenCeres: When Open Information Extraction Meets the Semi-Structured Web. *Proceedings of NAACL-HLT* (2019).

[46] Colin Lockard, Prashant Shiralkar, Xin Luna Dong, and Hannaneh Hajishirzi. 2020. ZeroShotCeres: Zero-Shot Relation Extraction from Semi-Structured Webpages. *ACL* (2020).

[47] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, and Sayak Paul. 2022. PEFT: State-of-the-art Parameter-Efficient Fine-Tuning methods. https://github.com/huggingface/peft.

[48] Mausam, Michael Schmitz, Stephen Soderland, Robert Bart, and Oren Etzioni. 2012. Open Language Learning for Information Extraction. *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning.*

[49] Avanika Narayan, Ines Chami, Laurel Orr, Simran Arora, and Christopher Ré. 2022. Can Foundation Models Wrangle Your Data? *Proceedings of the VLDB Endowment International Conference on Very Large Databases* (2022).

[50] Fatemeh Nargesian, Erkang Zhu, Reneé J. Miller, Ken Q. Pu, and Patricia C. Arocena. 2019. Data Lake Management: Challenges and Opportunities. *Proceedings of the VLDB Endowment* (2019).

[51] Christina Niklaus, Matthias Cetto, André Freitas, and Siegfried Handschuh. 2018. A Survey on Open Information Extraction. In *Proceedings of the 27th International Conference on Computational Linguistics.*

[52] OpenAI. March 2023. OpenAI API. https://openai.com/api/

[53] Laurel Orr. 2022. Manifest. https://github.com/HazyResearch/manifest.

[54] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155* (2022).

[55] F. Chen A. Doan P. DeRose, W. Shen and R. Ramakrishnan. 2007. Building structured web community portals: A top-down, compositional, and incremental approach. *VLDB* (2007).

[56] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. *arXiv:1606.05250* (2016).

[57] Alexander Ratner, Stephen H. Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré . 2017. Snorkel: Rapid Training Data Creation with Weak Supervision. *Proceedings of the VLDB Endowment (VLDB)* (2017).

[58] C. Romero and S. Ventura. 2013. Data mining in education. In *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*.

[59] Shreya Shankar, Rolando Garcia, Joseph M. Hellerstein, and Aditya G. Parameswaran. 2022. Operationalizing Machine Learning: An Interview Study. *arXiv:2209.09125* (2022).

[60] Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, Daniel Y Fu, Zhiqiang Xie, Beidi Chen, Clark Barrett, Joseph E Gonzalez, et al. 2023. High-throughput Generative Inference of Large Language Models with a Single GPU. *arXiv preprint arXiv:2303.06865* (2023).

[61] Jaeho Shin, Sen Wu, Feiran Wang, Christopher De Sa, Ce Zhang, and Christopher Ré. 2015. Incremental knowledge base construction using deepdive. In *Proceedings of the VLDB Endowment International Conference on Very Large Data Bases (VLDB)*.

[62] Ryan Smith, Jason A. Fries, Braden Hancock, and Stephen H. Bach. 2022. Language Models in the Loop: Incorporating Prompting into Weak Supervision. *arXiv:2205.02318v1* (2022).

[63] Immanuel Trummer. 2022. CodexDB: synthesizing code for query processing from natural language instructions using GPT-3 codex. *Proceedings of the VLDB Endowment* 11 (2022). https://doi.org/10.14778/3551793.3551841

[64] S. Raghavan S. Vaithyanathan T.S. Jayram, R. Krishnamurthy and H. Zhu. 2006. Avatar information extraction system. *IEEE Data Eng. Bull* (2006).

[65] Paroma Varma and Christopher Ré. 2018. Snuba: Automating Weak Supervision to Label Training Data.

[66] Paroma Varma, Frederic Sala, Ann He, Alexander Ratner, and Christopher Re. 2019. Learning Dependency Structures for Weak Supervision Models. *Proceedings of the 36th International Conference on Machine Learning (ICML)*.

[67] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le Le, Ed H. Cho, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-Consistency Improves Chain of Thought Reasoning in Language Models. *arXiv:2203.11171v2*.

[68] Eric Wu, Kevin Wu, Roxana Daneshjou, David Ouyang, Daniel Ho, and James Zou. 2021. How medical AI devices are evaluated: limitations and recommendations from an analysis of FDA approvals. *Nature Medicine* 27 (04 2021), 1–3.

[69] Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts. In *CHI Conference on Human Factors in Computing Systems*. 1–22.

[70] Wael M.S. Yafooz, Siti Z.Z. Abidin, Nasiroh Omar, and Zanariah Idrus. 2013. Managing unstructured data in relational databases. In *2013 IEEE Conference on Systems, Process & Control (ICSPC)*.

[71] Shaowen Zhou, Bowen Yu, Aixin Sun, Cheng Long, Jingyang Li, Haiyang Yu, Jian Sun, and Yongbin Li. 2022. A Survey on Neural Open Information Extraction: Current Status and Future Directions. *IJCAI22* (2022).

[72] Diana M. Zuckerman, Paul Brown, and Steven E. Nissen. 2011. Medical Device Recalls and the FDA Approval Process. *Archives of Internal Medicine* 171, 11 (06 2011), 1006–1011.

## A  EXPERIMENTAL DETAILS

### A.1  Evaluation Protocol

Because the baselines we compare against on SWDE require training data, they perform website-wise cross validation (*i.e.* train on some websites and evaluate on others). They do this for several combinations such that every website appears in the evaluation set. EVAPORATE, in contrast, does not require any training data, so we do not perform cross validation. To ensure a fair comparison, we simply evaluate EVAPORATE on all websites so that our method is evaluated on the same examples as the baselines.

### A.2  Metrics

We describe the metrics we use to evaluate OpenIE and ClosedIE performance of our system.

**Pair F1** For OpenIE, we report Pair F1 scores. Pair F1 is the standard metric for OpenIE systems [45]. The metric constructs (subject, value, predicate). The subject is the document-filename in our setting. The predicate is the attribute and the value is the attribute value. The F1 score computes the F1 score between the sets of gold and predicted tuples. This assigns credit for *exact-matches* between the attribute names and values extracted by the system and the ground truth — it assigns no partial credit.

Note that because EVAPORATE first identifies a list of attributes then sequentially generates functions and extracts the values, the user can "stop" execution at any number of attributes. The stopping point determines the number of tuples included in the prediction set. This is not a property of prior systems that extract "all or no" tuples [22, 39, 45, 46, 71, inter alia.]. For fair comparison we report performance at the number of gold attributes contained in the benchmark — note that this is generally *not* the number of attributes that maximizes the EVAPORATE's Pair F1 score.

**Text F1** For ClosedIE, we report Text F1 scores. Text F1 is the standard metric for extractive tasks and we use the exact implementation released by Rajpurkar et al. [56]. The metric tokenizes the prediction and gold strings and computes a token-wise F1 score.

Recall we select the F1 at the number of gold attributes, rather than at the number that gives the highest score).

## B  DATASET CONSTRUCTION

Below we describe how each of the evaluation benchmarks is obtained. We also release the suite of benchmarks along with the system code.

### B.1  FDA

For the FDA setting, we randomly sampled a dataset of 100 FDA 510(k) premarket notification submission PDFs for medical devices with substantially equivalent predicate devices since 1996 from the FDA website: https://www.accessdata.fda.gov/scripts/cdrh/cfdocs/cfpmn/pmn.cfm. We used the lightweight fitz library to convert this to text files. We asked 5 database graduate students to identify important attributes, defined as attributes useful for analysis that's present in at least a majority of documents. We collected the final set of 16 attributes as attributes agreed on by all graduate students. For these 16 attributes, we manually wrote functions to extract their value, then corrected errors by manual review. We defined

the attribute value as the full content of information pertaining to that attribute, typically a value, sentence, or section.

### B.2  Wiki NBA

For the Wiki NBA setting, we used the following SPARQL query over Wikidata to retrieve NBA articles. We then manually supplemented missing pages and filtered the results to only include pages about NBA players.

```
# Q13393265 is for Basketball Teams
# Q155223 is for NBA
# P118 is league (https://www.wikidata.org/wiki/
Property:P118)
SELECT ?item ?itemLabel ?linkcount WHERE {
    ?item wdt:P118 wd:Q155223 .
    ?item wikibase:sitelinks ?linkcount .
FILTER (?linkcount >= 1) .
SERVICE wikibase:label { bd:serviceParam wikibase
:language "[AUTO_LANGUAGE],en" . }
}
GROUP BY ?item ?itemLabel ?linkcount
ORDER BY DESC(?linkcount)
```

We asked 5 database graduate students to identify important attributes, defined as attributes useful for analysis that's present in at least a majority of documents. We collected the final set of 19 attributes as the attributes agreed on by all graduate students. For these 19 attributes, we manually wrote functions to extract their value, then corrected errors by manual review. We defined the attribute value as the full content of information pertaining to that attribute, typically a value, sentence, or section. We use these as ground truth extractions in the main paper.

We noted that the resulting attributes were complex for multiple documents, so we included another set of ground truth. We asked a graduate student to write functions to parse compound attributes whose values mention multiple values (e.g. birth date and location under attribute "Born") into atomic attributes and values. We this ground truth to demonstrate an additional step of schema cleaning in Section F.2.

### B.3  Enron

We download the Enron corpus from http://www.cs.cmu.edu/~enron/ and apply no further processing. We generate a benchmark using all metadata in the email headers by manually writing functions.

### B.4  SWDE

We download SWDE from https://github.com/cdlockard/expanded_swde. The benchmark includes the raw HTML from several websites with no further processing and ground-truth labels for selected attributes [45]. Because all the attributes are located within the root-elements of the HTML *body*, excluding the information such as the HTML header, attributes described within tags (e.g. <a href='/year/2012/>', <title>), and so on, we extend the original benchmark to include a more diverse set of attributes. We refer to the extended benchmark as SWDE Plus.

## C WEAK SUPERVISION DETAILS

*Objective.* We detail the weak supervision (WS) algorithm used for aggregating the votes across generated functions. Let $\mathcal{D}$ be our *unlabeled* dataset of documents from which we are extracting a particular attribute of interest. Let $y$ be a random variable representing the true attribute value. Let $\lambda$ represent the outputs of our $m$ generated extraction functions $f \in \mathcal{E}$ on a particular document. Each $\lambda_i \in \lambda$ is a function $\lambda_i : \mathcal{X} \rightarrow \mathcal{Y}$. Our goal is to use the vectors $\lambda$, produced across documents $x \in \mathcal{D}$ to infer the true label $y$. Concretely, we seek, $\phi(x)$, a function that inputs $\lambda$ and outputs a final prediction $\hat{y}$ for a document $x$.

With labeled data, i.e. where we had documents and their attribute values $\{(x_1, a_1), ..., (x_n, a_n)\}$, we could perform traditional supervised learn $\phi(x)$ that maps $\lambda$ to $y$. However, in our unlabeled setting, the insight is to use the noisy estimates of $y$ produced by each of the functions to construct $\phi(x)$.

*Standard WS Setup [57].* WS models learn the latent variable graphical model on the distribution $\Pr(y, \{\lambda_1, ..., \lambda_m\} = \Pr(y, \lambda)$. In this model, $y$ is *latent*, we cannot observe it. To produce $\hat{y}$, we concretely perform two steps:

- **Learn the label model** We have access to $\lambda$ and can use this to learn the *label model* $P(\lambda|y)$.
- **Inference** We can then produce the estimate $\hat{y}$ by setting $\phi(x) = \text{argmax}_y \Pr(y|\lambda(x)$

First, we will discuss how to learn the label model. We parameterize $P(\lambda|y)$ as follows:

$$P(\lambda|y) = \frac{1}{Z_\theta} \exp(\sum_{i=1}^{m} \theta d(\lambda_i, y))$$

Intuitively when modeling our sources, we want to model how accurate a particular $\lambda_i$ is, when it makes a prediction. We are then going to want to weigh the predictions of the different sources proportional to their accuracies.

$Z$ is a constant to normalize the probability distribution. The feature-vector $d$ can be broken down into:

- $d_{\text{lab}}$, representing how frequently $\lambda_i$ provides a prediction vs. abstains across documents. This **is** directly observable.
- $d_{\text{corr}}$, which represents how frequently $\lambda_i$ and $\lambda_j$ yield the same prediction for the attribute value. This **is** directly observable.
- $d_{\text{acc}}$, representing the accuracy of $\lambda_i$, or how frequently $\lambda_i$ agrees with $y$. Note that the accuracy is measured across documents for which the function provides a prediction and does not abstain. This **is not** directly observable.

$\theta$ are the parameters we aim to learn to combine the inputs in the feature vector. To learn this without access to $y$, we can minimize the negative log marginal likelihood given the observed $\lambda_i$ outputs and solve with SGD, or use a closed-form solution [27].

### C.1 Theoretical Analysis

Here we analyze our proposed weak supervision algorithm (Algorithm 1) for EVAPORATE-CODE+.

*Setup.* Recall that in Algorithm 1, the LLM is 1) used to generate the candidate extraction functions $F = \{f_1, ..., f_m\}$, and 2) to directly extract values from a small set of $n$ documents, $D_{eval}$. Let $\{\hat{y}_1, ..., \hat{y}_n\}$ be the LLM's direct extractions over $D_{eval}$. The candidate functions are each applied to $x_i \in D_{eval}$ to produce $\{f_i(x_1), ..., f_i(x_n)\}$ for each $f_i \in F$.

The candidate function outputs $\{f_i(x_1), ..., f_i(x_n)\}$ are then scored against the LLM's direct extractions $\{\hat{y}_1, ..., \hat{y}_n\}$ to estimate the function accuracy:

$$\hat{a} = \frac{1}{n} \sum_{j=1}^{n} \mathbf{1}(f_i(x_j) = \hat{y}_j)$$

From Varma and Ré [65], we can estimate whether the accuracies learned via weak supervision $\tilde{a}$ will be close to the true function accuracies $a^*$. We can use the function accuracy estimated on $D_{eval}$ to estimate whether the accuracies $\tilde{a}$ learned by the weak supervision model $\theta$ are close to the true function accuracies $a^*$, i.e. we can bound $||a^* - \tilde{a}||_\infty < \gamma$. Concretely, the objective is to provide a guarantee that $\gamma$ is the upper bound on the true error, given that the measured error $||\hat{a} - \tilde{a}||_\infty$ is below some threshold $\epsilon$. Proposition 1 in Varma and Ré [65] states the following. Refer to Varma and Ré [65] for the proof.

*We have $m$ functions with empirical accuracies $\hat{a}$ and learned accuracies $\tilde{a}$, and the measured error is below some threshold $\epsilon$. Then, if each function labels a minimum of*

$$n \geq \frac{1}{2(\gamma - \epsilon)} \log(\frac{2m^2}{\delta})$$

*datapoints, the weak supervision label model will succeed in learning accuracies such that $||a^* - \tilde{a}||_\infty < \gamma$ with a probability $1 - \delta$.*

*Extension.* However, in our setting, $\hat{y}_i$ may not be equal to the true extraction $y_i^*$ for $d_i \in D_{eval}$ because the LLM is being used to score the candidate functions, rather than a gold-standard labeled dataset. Suppose the LLM introduces its own error rate $e$, then our objective is to modify the above statement to account for this source of error. Intuitively, as $e$ increases, we need more datapoints to obtain a better estimate of the accuracies.

*Assumptions and Notation* We will assume the errors are randomly occurring across $x_i \in D_{eval}$ and that the data points in $D_{eval}$ are independent. Let $\hat{a}$ be the accuracy estimated using the noisy labels $\hat{y}$ on $D_{eval}$ and let $a^\star$ be the function accuracy given ground truth labels $y$ on $D_{eval}$. We follow from Varma and Ré [65] below, but extend the theory to support our setting.

*Proposition 1. We have $m$ functions with empirical accuracies $\hat{a}$, where the gold labels have error rate $e$, the learned accuracies are $\tilde{a}$, and the measured error is below some threshold $\epsilon$. Then, if each function labels a minimum of*

$$n \geq \frac{1}{2(\gamma - \epsilon - e)} \log(\frac{2m}{\delta})$$

*datapoints, the weak supervision label model will succeed in learning accuracies such that $||a^* - \tilde{a}||_\infty < \gamma$ with a probability $1 - \delta$.*

*Proof.* We can decompose the true error $||a^* - \tilde{a}||_\infty$ using the triangle inequality:

$$P(||a^* - \tilde{a}||_\infty > \gamma) \leq P(||a^* - \hat{a}||_\infty + ||\hat{a} - \tilde{a}||_\infty > \gamma)$$
$$\leq P(||a^* - \hat{a}||_\infty + \epsilon)$$

Recall that for candidate function $f_i$, where $n_i$ is the number of data points on which $f_i$ did not abstain:

$$\hat{a}_i = \frac{1}{n_i} \sum_{j:f_i(x_j) \neq \emptyset} \mathbf{1}(f_i(x_j) = \hat{y}_j)$$

where $\hat{y}_j$ is the direct LLM extraction on document $x_j \in D_{eval}$, $f_i(x_j)$ is the function output on $x_j$, and $n_i$ is the set of examples where the function did not abstain (i.e. $f_i(x_j) \neq \emptyset$). Suppose $P(\hat{y}_j \neq y_j^*) = e$. We again apply the triangle inequality to include $e$ in our decomposition of the true error.

$$P(||a^* - \tilde{a}||_\infty > \gamma) \leq P(||a^* - a^\star||_\infty + ||a^\star - \hat{a}||_\infty + ||\hat{a} - \tilde{a}||_\infty > \gamma)$$

$$\leq P(||a^* - a^\star||_\infty + \epsilon + e)$$

Taking a union bound across the $m$ functions:

$$P(||a^* - \hat{a}||_\infty + \epsilon + e > \gamma) \leq \sum_{i=1}^{m} P(||a_i^* - a^\star||_\infty + \epsilon + e > \gamma)$$

Next, $a^\star$ for function $f_i$ is defined, with respect to the ground truth labels $y$ on $D_{eval}$, where $n_i$ is the number of data points on which $f_i$ did not abstain:

$$a^\star = \frac{1}{n_i} \sum_{j:f_i(x_j) \neq \emptyset} \mathbf{1}(f_i(x_j) = y_j)$$

Substituting $a^\star$ into our bound:

$$P(||a^* - \hat{a}||_\infty + \epsilon + e > \gamma) \leq \sum_{i=1}^{m} P(||a_i^* - a^\star||_\infty + \epsilon + e > \gamma)$$

$$\leq \sum_{i=1}^{m} P(|a_i^* - \frac{1}{n_i} \sum_{j:f_i(x_j) \neq \emptyset} \mathbf{1}(f_i(x_j) = y_j)| > \gamma - +\epsilon - e)$$

Next we can apply Hoeffding's inequality:

$$\leq \sum_{i=1}^{m} 2 \exp(-2(\gamma - \epsilon - e)^2 n_i)$$

$$\leq 2m \exp(-2(\gamma - \epsilon - e)^2 \min(n_1, ... n_m))$$

Finally, rearranging the terms to isolate $n$, where $n = \min(n_1, ... n_m)$:

$$P(||a^* - \tilde{a}||_\infty > \gamma) \leq 2m \exp(-2(\gamma - \epsilon - e)^2 n)$$

$$\delta \leq 2m \exp(-2(\gamma - \epsilon - e)^2 n)$$

Taking the log of both sides each side is $< 1$ in value:

$$\log(\delta) \geq \log(2m \exp(-2(\gamma - \epsilon - e)^2 n))$$

$$\log(\delta) \geq \log(2m) + \log(\exp(-2(\gamma - \epsilon - e)^2 n))$$

Negating both sides:

$$\log(2m) - \log(\delta) \leq 2(\gamma - \epsilon - e)^2 n$$

$$\frac{1}{2(\gamma - \epsilon - e)^2} \log(\frac{2m}{\delta}) \leq n$$

*Analysis.* As the number of aggregated candidate functions $m$ increases, then the number of points each function should label for learning the label model, i.e. the size of $D_{eval}$, should be increased. As the underlying LLM generating $\hat{y}$ increases in accuracy, we can decrease the size of $D_{eval}$.

**Table 6: Quality and cost achieved through prompting OpenAI's Text-Davinci-003 model to extract specific, pre-defined attributes.**

| Source (Format) | Quality | | Cost / 10k Documents | |
|---|---|---|---|---|
| | # Attributes | F1 | Tokens (M) | Dollars ($) |
| Enron Emails (TXT) | 15 | 85.3 | 140 | 2,790 |
| FDA (TXT) | 16 | 78.0 | 241 | 4,816 |
| Wiki NBA (HTML) | 19 | 84.6 | 328 | 6,559 |
| SWDE Movies (HTML) | 25 | 84.4 | 359 | 7,174 |
| SWDE University (HTML) | 33 | 72.6 | 379 | 7,586 |
| **Average** | **21.6** | **79.9** | **289** | **5,785** |

**Table 7: The left column includes our main results from Table 1 for reference. On the right, we evaluate the quality of EVAPORATE when ground-truth labels are used instead to score functions and estimate the abstention fraction $e$ in Algorithm 1.**

| Source (Format) | OpenIE with $\mathcal{F}$ | OpenIE with Ground-Truth |
|---|---|---|
| Enron Emails (TXT) | 89.9 | 87.0 |
| FDA (TXT) | 62.8 | 61.9 |
| Wiki NBA (HTML) | 68.2 | 79.7 |
| SWDE Movie (HTML) | 56.8 | 64.0 |
| SWDE University (HTML) | 59.0 | 66.3 |
| **Average** | **66.7** | **71.8** |

# D ADDITIONAL ABLATIONS FOR EVAPORATE

## D.1 Validating the Quality of FM $\mathcal{F}$

Because our approach for scoring the generated functions relies on comparing to the extractions produced by directly processing the document with the FM $\mathcal{F}$ (Section 3.1), in Table 6, we evaluate the quality of $\mathcal{F}$. In contrast to the "End-to-End" prompt baseline (Prompt ??), in this experiment, the FM is prompted with a specific attribute to extract. This prompt is shown in Appendix ??.

Next, we compare the quality that is achieve using Algorithm 1 by using *ground truth labels* on $D_{eval}$ (containing 10 documents) to EVAPORATE achieves using FM $\mathcal{F}$. Recall that in Algorithm 1, the "labels" from $\mathcal{F}$ are used to both estimate the fraction $e$ and the scores for functions $f \in F$. We find that the upper bound quality, using ground-truth labels, is 5.1 F1 points higher as shown in Table 7.

## D.2 Varying Number of Candidate Functions

In this section, we explore how varying the number of candidate functions generated by EVAPORATE-CODE+ affects the quality and efficiency of the system.

Recall from 3.3.2, that EVAPORATE-CODE+ generates a diverse set of candidate extraction functions a diverse set of *function candidates* $F = \{f_1, f_2, ... f_k\}$, with variable quality (Figure 6). The outputs of the generated functions are then aggregated using a weak supervision algorithm (Algorithm 1).

In our main experiments, we generate the pool of candidate functions and score the candidate functions on $d = 10$ documents.

**Figure 6: Across the gold attributes and settings, we report the estimated quality of the generated extraction functions.**

**Table 8: Quality using varying numbers of documents from which candidate functions are generated and scored.**

| Source Dataset (Format) | 1 | 3 | 5 | 10 |
|---|---|---|---|---|
| FDA (TXT) | 40.0 | 56.5 | 67.1 | 62.8 |
| Enron Emails (TXT) | 45.5 | 54.1 | 67.2 | 86.9 |
| Wiki NBA (HTML) | 64.3 | 61.4 | 64.7 | 68.2 |
| SWDE Movies (HTML) | 45.6 | 46.7 | 51.7 | 56.8 |
| SWDE Universities (HTML) | 38.9 | 44.1 | 55.1 | 59.0 |
| **Average** | **46.9** | **52.6** | **61.2** | **66.7** |

**Table 9: Quality achieved by passing varying numbers $k$ of the top-$k$ candidate functions to the weak supervision algorithm. For all runs, the pool of candidate functions are generated and scored by prompting on 10 seed documents.**

| Source Dataset (Format) | 1 | 3 | 5 | 10 |
|---|---|---|---|---|
| FDA (TXT) | 49.0 | 59.1 | 62.1 | 62.8 |
| Enron Emails (TXT) | 85.7 | 85.7 | 85.7 | 86.9 |
| Wiki NBA (HTML) | 69.9 | 67.2 | 68.9 | 68.2 |
| SWDE Movies (HTML) | 56.9 | 55.7 | 55.4 | 56.8 |
| SWDE Universities (HTML) | 56.7 | 58.9 | 59.8 | 59.0 |
| **Average** | **63.6** | **65.3** | **66.4** | **66.7** |

We then feed the top $k = 10$ highest scoring candidate functions to the weak supervision algorithm.

First, we measure the quality of EVAPORATE-CODE+ as we vary the number of documents $d$ used to generate and score the candidate functions. We evaluate for $d \in \{1, 3, 5, 10\}$. The results are reported in Table 8.

Next, we measure the quality of EVAPORATE-CODE+ as we vary the number of functions $k$ passed to the weak supervision algorithm. We evaluate for $k \in \{1, 3, 5, 10\}$. The results are reported in Table 9.

The quality increases with both the number of seed documents used to generate and evaluate candidate functions, and with the number of candidate functions passed to the weak supervision algorithm.

# E  ADDITIONAL COMPARISONS TO BASELINES

Here we provide additional quality and efficiency comparisons between EVAPORATE and prior work.

## E.1  Additional Baselines from Prior Work

Here we study two additional baselines from prior work for OpenIE and ClosedIE respectively, beyond the Document-level IE methods discussed in Section 4.3 (Table 2).

*OpenIE.* We apply the OpenIE6 system from Kolluru et al. [39], which is a state-of-the art approach for OpenIE over unstructured (non-HTML) text. While this system is not designed for extracting structured views over heterogeneous data, we evaluate it qualitatively to understand how existing OpenIE systems perform in this setting.

(1) First, the system only handles well formatted sentences and struggles to extend to heterogeneous data types. It does not handle HTML documents and is difficult to apply to documents with complex formatting (like PDFs) where full sentences can be difficult to extract. Using SWDE College Prowler as an example, given the HTML input line `<td>Student Body Size:</td> <td class="stat"> <span id="..."> 6,504 </span> </td>`, OpenIE6 misses the student body size attribute and corresponding value.

(2) Second, even when documents contain full sentences, OpenIE6 extracts an extremely large set of relations for each document and does not prioritize attributes by relevance or enforce consistent attributes across documents. This makes it difficult to use the resulting relations to understand the contents of the documents or to do analysis. Using a sample FDA 510(k) document from our corpus as an example, OpenIE6 extracts 427 relations with 184 relations with confidence larger than 0.5. While some can be informative, a significant fraction of these relations are not useful for indexing and analysis across documents. For example, "(sample carryover; occurs; the results of the acceptor assay show interference)" is an extracted relation with confidence 0.99.

*ClosedIE.* We study the effectiveness of span-extractor models, which are commonly used in QA systems to extract the information that is relevant to a user-query from a provided document context [19]. Given the ground truth attributes, we evaluate the ability of these models to extract their values from the relevant paragraphs. We evaluate the DebertaV3 Large model fine-tuned on the Squad 2.0 dataset, which achieves 90.8 F1 on the Squad 2.0 dev set in Table 10. We find text-davinci-003 (Table 6) and our EVAPORATE function generation approach (Table 1) significantly outperforms this pre-trained QA model on ClosedIE in all settings, over text and HTML documents.

**Table 10: Results as in Table 6 using the DebertaV3 large model fine-tuned on the Squad2.0 dataset from HuggingFace.**

| Source (Format) | # Attributes | Closed IE F1 |
|---|---|---|
| Enron Emails (TXT) | 15 | 53.7 |
| FDA (TXT) | 17 | 56.5 |
| Wiki NBA (HTML) | 19 | 50.2 |
| SWDE Movies (HTML) | 30 | 43.5 |
| SWDE University (HTML) | 25 | 45.3 |

## E.2 Efficiency Comparisons Between EVAPORATE and Baselines

In this section, we compare the efficiency of EVAPORATE and baselines like DOM-LM [22]. First, we breakdown the efficiency analysis in terms of the parameter count, pre-training cost, training costs, and inference costs. Then we use the cost break down to quanitatively compare EVAPORATE to the baselines Table 11. We also quantitatively compare EVAPORATE-DIRECT and EVAPORATE-CODE+, highlighting how the decision hinges on the number of documents to be processed, as well as the number of attributes in the dataset. Finally, we compute estimated FLOP counts required to process documents from the five sources we considered (see Table 12). We use GPT3-175B as the language model in our calculations of costs for EVAPORATE, since those are the primary results we report in the paper. We report experiments with multiple models in the paper and note that EVAPORATE is flexible with respect to the underlying language model. This allows users to tradeoff quality and costs by using smaller language models (e.g. the smallest GPT model has with 175 million parameters) [28].

- **Parameter count** The baseline uses a RoBERTa-Base LM backbone model, which is 125M parameters, while EVAPO-RATE uses a 175B model.[1] We note that the baselines require training a new model for each new data domain, so the final parameter count for the baselines is (Number of Domains× Model Size).
- **Pre-training costs** The LLMs used in EVAPORATE see fewer tokens (300 billion tokens) [14, 54] than RoBERTa-base (2 trillion tokens) during pretraining [44]. Brown et al. [14] compares the pretraining of both models. Because of the disparity in parameter counts discussed above, the pre-training phase of the models used in EVAPORATE is roughly 200-times more FLOP-intensive than the baseline. However, users of EVAPORATE or the baselines are not responsible for pre-training. Pre-trained models are commonly available via model hubs (*e.g.* HuggingFace) and APIs, so pre-trained models can be reused in a wide variety of applications and by many users, amortizing the cost.
- **Training Costs** EVAPORATE does not require any task-specific training while the baselines require that users train a new model for each data domain. Unlike pretraining, the task-specific training of the baselines would likely need to be performed on users' local compute. The exact cost

of training is hard to estimate as it depends on the number of training epochs and the amount of training data. Additionally, collecting the requisite training data is a manual and time-intensive process. Specifically, SWDE Movies has data from 8 websites and SWDE Universities has data from 5 websites. Under the "zero-shot" protocol specified in Deng et al. [22], DOM-LM trains on 10% of the data from 6 and 3 websites respectively for Movies and Universities, and evaluates on the heldout websites within each domain. Nonetheless, DOM-LM underperforms EVAPORATE, which uses no task-specific labeled training data.

- **Inference costs** Inference FLOP requirements grow linearly in the model parameter count, so the cost of running inference on a single document with an LLM is significantly more expensive than with a baseline model. As shown in Table 11, both DOM-LM and EVAPORATE-DIRECT require $O(n)$ inference calls to process $n$ documents. On the other hand, EVAPORATE-CODE+ only requires $O(m)$ inference calls in order to generate extraction functions which can then be applied to the entire corpus of documents. Thus, as the the number of documents to be processed grows, the inference cost of baseline methods eventually surpass the inference costs of EVAPORATE-DIRECT.

*Quantifying the costs.* First we note that the FLOP requirement for pretraining RoBERTa and GPT-3 are reported in [14]:

- **RoBERTa** The total training FLOPS is 1.50E+21. The inference FLOPS per token is 2 × Number of Parameters, which is 0.250 GFLOPS.
- **GPT3-175B** The total training FLOPS is 3.14E+23. The inference FLOPS per token is 2 × Number of Parameters, which is 350 GFLOPS.

The cost of applying EVAPORATE and baselines to a document corpus varies depending on the number of documents $n$, the number of attributes $m$, and the model architecture used. The costs are summarized in Table 11.

- **DOM-LM** The inference cost is computed by

$$n \times \frac{\text{tokens}}{\text{document}} \times 0.250 \text{GFLOP}$$

where 0.250GFLOP is the inference cost per token of DOM-LM.
- **EVAPORATE-CODE+** The inference cost is computed by

$$m \times P \times \frac{\text{tokens}}{\text{chunk}} \times 350 \text{GFLOP}$$

where $P$ is the number of inference calls required per attribute and 350 GFLOP is the inference cost per token of EVAPORATE, see Table 11. Note that for our evaluated implementation $P \approx 10c$ for a small number $c$. To understand where $c$ comes from, recall $D_{eval}$ is 10 documents. We generate candidate functions on chunks of each document and directly apply the LLM to each document so that we can score the function outputs against the LLM outputs.).

The inference FLOP counts for DOM-LM and EVAPORATE-CODE+ fall in the same order of magnitude, though which method is more

---

[1]While OpenAI has not released models, the `text-davinci-003` model is estimated to be 175B parameters [28].

**Table 11: Statistics on computational cost of the language models used by EVAPORATE and baselines. The parameter counts and pre-training costs are sourced from [14]. *Inference calls* is the number of LLM forward passes required to process $n$ documents with $m$ attributes. FLOP counts assume that DOM-LM is implemented with a RoBERTa backbone (as reported in the paper [22]) and that EVAPORATE is implemented with text-davinci-003.**

| | Parameter Count (B) | Pre-training Cost (ZFLOP) | Fine-tuning Cost (GFLOP per token) | Inference Calls (# of Documents) | Inference Cost (GFLOP per token) |
|---|---|---|---|---|---|
| DOM-LM [22] | 0.125 | 1.50 | 11.250 | $O(n)$ | 0.250 |
| EVAPORATE-DIRECT | 175 | 341 | 0 | $O(n)$ | 350 |
| EVAPORATE-CODE+ | 175 | 341 | 0 | $O(m)$ | 350 |

**Table 12: Language model inference costs for processing 100K Documents with EVAPORATE and baselines. Because document length and number of attributes, $m$, varies by source, we report costs for each of the sources included in our experiments. FLOP counts are based on the statistics reported in Table 11. FLOP counts assume that DOM-LM is implemented with a RoBERTa backbone (as reported in the paper [22]) and that EVAPORATE is implemented with text-davinci-003.**

| | Cost / 100K Documents | | | | | |
|---|---|---|---|---|---|---|
| Source (Format) | DOM-LM [22] | | EVAPORATE-DIRECT | | EVAPORATE-CODE+ | |
| | Tokens (M) | PFLOP | Tokens(M) | PFLOP | Tokens (M) | PFLOP |
| FDA (TXT) | 145.6 | 364 | 145.6 | 50,960 | 1.9 | 665 |
| Enron Emails (TXT) | 21.2 | 53 | 21.2 | 7,420 | 0.6 | 210 |
| WIki NBA (HTML) | 650.1 | 1,625 | 650.1 | 227,535 | 3 | 1,005 |
| SWDE Movie (HTML) | 282.9 | 707 | 282.9 | 99,015 | 2.3 | 805 |
| SWDE University (HTML) | 190.1 | 475 | 190.1 | 66,535 | 1.9 | 665 |

costly depends on the document source and the number of attributes. Users can select the most efficient method for their setting using the above cost framework.

## F  EXTENSIONS TO EVAPORATE

In this section, we provide experiments for three extensions to the EVAPORATE system including 1) operating over a dataset that contains a mixture of the document sources (e.g. a mixture of SWDE Movies websites) (Section F.1), 2) further cleaning the output extractions from EVAPORATE (for instance by converting them into atomic attributes) (Section F.2) and 3), using domain specific prompts to further control the system quality (Section F.3). and a summarizing discussion of future directions F.4.

### F.1  Recovering Document Sources

In this section, we demonstrate how we can recover the sources from an unlabeled mixture of documents.

In our experiments, we process each source with EVAPORATE separately, instead of mixing them together. This more closely matches the experimental setup of the baselines like DOM-LM [22]. However, in practice, we may not have access to the sources of the documents in our data lake. To verify that we would be able to recover document sources from an unlabeled data lake, we explore whether we can use unsupervised methods to group documents by source.

First, we mix the documents from the websites of the Movie and University subsets of SWDE. This provides an unlabeled "data lake" with distinct document sources. Note that many of the sources come from the same domain, towards increasing the difficulty of

distinguishing the sources. We apply a standard TF-IDF vectorizer and K-means clustering to the mixture of documents. Finally, we measure the adjusted rand index (ARI), a standard clustering metric which allows for random permutations in the cluster labels. We achieve a perfect ARI of 1.0, meaning we were able to exactly recover the document sources without any labeled data or supervision.

In Figure 5, we include a T-SNE visualization of the document TF-IDF vectors. This plot illustrates the separation in representation space between documents from different websites.

### F.2  Atomic Schema Cleaning Extensions

One further extension of schema generation is atomic schema cleaning. EVAPORATE generates a set of candidate attributes which, in some cases, are complex and can be decomposed into cleaned, atomic attributes. For example, an extracted attribute of "born" from the Wiki NBA setting, has the following form: <birth_date> (age) <location>. Decomposing the "born" attribute into three separate attributes (e.g., birth date, age and birth location) would enable users to ask queries such as — How many players in the NBA were born in Ohio? — that would otherwise be unanswerable with the existing schema. As such, decomposing the complex attributes into cleaned, atomic attributes increases the utility of the resulting schema and extractions for analysis and indexing. Prior work [49] has demonstrated that FMs can be useful for data transformation task. Schema decomposition can be viewed as an instantiation of data transformation, suggesting that such an operation could be completed using a FM.

We manually clean the ground truth complex attributes and values in the Wiki NBA setting and construct the ground truth

atomic attribute and values. We find that after cleaning there are 35 atomic attributes for Wiki NBA, decomposed from the 19 complex attributes.

For our method, we prompt the expensive FM (in this case text-davinci-003 from OpenAI) to decompose the complex attribute and values into a list of atomic attributes and values, for a single example of each complex attribute and value. To save computation cost, we then use the large LLM schema cleaning result from one example to prompt a smaller, less expensive FM (in this case the text-curie-001 model from OpenAI) to extract the cleaned values for the remainder of documents. We provide the smaller, less expensive FM with the complex attribute, the cleaned attribute to extract, and a one-shot example from the expensive FM.

To measure the performance of schema cleaning, we construct pairs of (file, value) for all files and values and compute the precision, recall, and F1 as in the OpenIE setting against the ground truth atomic attributes and values. We do not include the attribute in the relations to score, because the extracted values are generally unique and we want to avoid penalizing generated atomic attribute names that differ from the ground truth but are still correct. As a baseline before our atomic schema cleaning step, we score the ground truth complex values against the ground truth atomic values and find it achieves an F1 of 21.0, since many values are not atomic. Applying our atomic schema cleaning methodology to the ground truth complex values decomposes them into atomic attributes, qualitatively improving the usability of the attributes. The resulting predicted atomic attributes achieve an F1 of 57.5 when scored against the ground truth atomic values.

### F.3 Improving Quality via Domain Specific Prompts

A potential direction to improve the quality of EVAPORATE is to replace the fixed, generic prompts with domain (e.g. Movies) or dataset (e.g. Rotten Tomatoes website) prompts.

Towards demonstrating this opportunity, we update the Schema Identification prompt (Prompt ??) to include in-context examples that discuss the Movies domain. The resulting prompt is provided in Section ??. We apply the single prompt to all 8 SWDE Movies websites (i.e. the prompts are not customized on a per-website level, but rather at the Movies domain-level). We observe this results in 4.5 Recall@K (8%) improvement, averaged across the 8 settings. Results by website are provided in Table 13.

### F.4 Discussion of Future Directions

The goal of this study was to evaluate a simple, prototype system that uses LLMs to generate structured views from unstructured documents. We explored two fundamentally different implementation strategies, highlighting opportunities for future work in the space.

**New applications** Our findings demonstrate the promise of function synthesis as a way to mitigate cost when using LLMs. We study the problem of materializing a structured view of an unstructured dataset, but this insight may be applicable in a broader suite of data wrangling tasks. Many data wrangling tasks are high throughput applications, for which LLMs are not optimized. Future

**Table 13: Comparing Recall@$K$, where $K$ is the number of gold attributes for the dataset, on Schema Identification using generic, domain agnostic vs. domain specific in-context demonstrations for the 8 SWDE Movie websites.**

| SWDE Movies Website | Base | Domain | Difference |
|---|---|---|---|
| Allmovie | 62.7 | 65.2 | + 2.5 |
| AMCTV | 88.2 | 83.9 | − 4.3 |
| Hollywood | 14.4 | 7.9 | − 6.5 |
| iHeart Movies | 62.5 | 75.0 | + 12.5 |
| IMDB | 70.4 | 71.6 | + 1.2 |
| Meta Critic | 35.3 | 64.7 | + 29.4 |
| Rotten Tomatoes | 66.7 | 58.3 | − 8.4 |
| Yahoo Movies | 72.7 | 81.8 | + 9.1 |
| **Average** | **59.1** | **63.6** | **4.5** |

work should explore whether *code synthesis* may enable general low cost solutions.

**New function types** We dichotomized the *direct extraction* and *code synthesis* implementations. However, the lines between these approach may blur going forward. After all, the LLM could generate functions that invoke *other models* — for instance, the LLM may generate functions that call the NLTK, HUGGINGFACE, or even the OPENAI APIs. This naturally raises the question of how to characterize the cost of the generated functions, rather than assuming they are inexpensive.

**Improving quality** Future work may consider iterative approaches to function generation. Concretely, when a generated function fails to compile or achieves low scores compared to the high quality LLM, we may be able to provide the compilation errors and/or high quality LLM responses in a prompt that encourages the LLM to generate an improved function. For instance, we may use Algorithm 1 to score the quality of *small* LLMs for performing the extractions.

The goal of our work is to provide improved methods and infrastructure for using recent LLM technology. A complementary objective is to improve the quality of the LLM itself. The LLMs with high quality in-context learning results tend to be parameter intensive, which means that further training the model can be expensive. There are two main approaches in the literature that explore efficiently adapting the LLM itself to new data that we refer to: parameter-efficient fine-tuning and retrieval augmented generation. Parameter-efficient fine-tuning methods update a small number of model parameters, keeping the rest frozen, thereby significantly improving the efficiency of fine-tuning while still attaining strong results [47]. Retrieval-augmented generation methods retrieve relevant information and incorporate it into generation, thereby making it easier to update and change the information used by the LLM [41].

## G EXAMPLE GENERATED CANDIDATE FUNCTIONS

### G.1 Case Study of Function Generation and Aggregation

Here we provide two end-to-end examples of 1) generated candidate functions, 2) the outputs collected from different candidate functions and fed into the weak supervision algorithm, and 3) the resulting outputs. The first is an example attribute that achieves high quality under EVAPORATE and the second is an attribute that achieves low quality.

**FDA Setting** First we consider the 510k attribute in the FDA reports dataset.

Five of the ten generated functions are shown below:

```
def get_510_k_number_field(text: str):
    """
    Function to extract 510(k) number.
    """
    pattern = r"K\d+"
    return re.findall(pattern, text)

def get_510_k_number_field(text: str):
    """
    Function to extract the 510(k) number.
    """
    pattern = r'510\(k\) Number:\s*(K\d+)'
    return re.findall(pattern, text)[0]

def get_510_k_number_field(text: str):
    """
    Function to extract 510(k) number.
    """
    pattern = r'(?:510\(k\))\s*number\(s\):\s
     *(.*)'
    result = re.findall(pattern, text)
    return result

def get_510_k_number_field(text: str):
    """
    Function to extract the 510(k) number.
    """
    pattern = r'K\d{6}'
    return re.findall(pattern, text)

def get_510_k_number_field(text: str):
    """
    Function to extract the 510(k) number.
    """
    pattern = r'510\(k\) Number:\s*(k\d+)'
    return re.findall(pattern, text)[0]
```

Five examples of relevant snippets from FDA reports in the dataset are shown below as examples for the case study:

```
Report: K143467.txt
Content snippet:
ASSAY AND INSTRUMENT COMBINATION TEMPLATE
A. 510(k) Number:
k143467


Report: K171742.txt
Content snippet:
DECISION MEMORANDUM

A. 510(k) Number:


K171742


Report: K170491.txt
Content snippet:
ASSAY AND INSTRUMENT COMBINATION TEMPLATE
A. 510(k) Number:
K170491


Report: K171641.txt
Content snippet:
DECISION SUMMARY
A. 510(k) Number:
K171641


Report: K162526.txt
Content snippet:
ASSAY ONLY TEMPLATE
A. 510(k) Number:
k162526
```

The resulting outputs of the functions applied to the reports are shown below. Note how 6 of 10 functions fail to extract the 510k number for report K143467.txt – simply taking the majority vote across the functions would yield an empty-value in the EVAPORATE output. Further, how one function provides an incorrect prediction "K2" on report K162526.txt. The goal of the weak supervision algorithm is to model the function predictions across the dataset to determine which to rely on for different documents.

```
Report: K143467.txt
Predictions: ['', '', '', '', '', '', 'k143467',
    'k143467', 'k143467', 'k143467']

Report: K171742.txt
Predictions: ['K171742', 'K171742', 'K171742', '
    K171742', 'K171742', 'K171742', '', '', '',
    '']

Report: K170491.txt
Predictions: ['K170491', 'K170491', 'K170491', '
    K170491', 'K170491', 'K170491', '', '', '',
    '']

Report: K171641.txt
Predictions: ['K171641', 'K171641', 'K171641', '
    K171641', 'K171641', 'K171641', '', '', '',
    '']

Report: K162526.txt
Predictions: ['', '', '', '', '', 'K2', 'k162526
    ', 'k162526', 'k162526', 'k162526']
```

We then apply the weak supervision algorithm (Algorithm 1) to the matrix of function outputs across all documents in the dataset (i.e. a matrix of shape $100 \times 10$ for the 100 documents and 10 candidate functions for the FDA setting). Because each function can output widely different values and the values differ across documents, the set of unique predictions per document varies. However, the label model expects a 1) fixed number of 2) constant (i.e. same for every document) classes. We handle the expectation for constant classes by *symmetrizing* the label model as described in Section 3.3. We handle the expectation for a fixed number of classes by taking the most frequently occurring $k$ (e.g. $k = 5$) classes per document. If there are $< k$ unique predictions for a documenet, we fill the input vector with dummy values. For instance, this results in the following, given the 10 outputs per document shown in the prior step:

```
Report: K143467.txt
Inputs: ['', 'k143467', 'dummy -1', 'dummy -2', '
    dummy -3']

Report: K171742.txt
Inputs: ['K171742', '', 'dummy -1', 'dummy -2', '
    dummy -3']

Report: K170491.txt
Inputs: ['K170491', '', 'dummy -1', 'dummy -2', '
    dummy -3']

Report: K171641.txt
Inputs: ['K171641', '', 'dummy -1', 'dummy -2', '
    dummy -3']

Report: K162526.txt
Inputs: ['', 'K2', 'k162526', 'dummy -1', 'dummy
    -2']
```

Finally, we learn the label model on the unlabeled inputs as shown above. The label model results in the following predictions on the five samples. The predictions are returned to the user in the output of EVAPORATE.

```
Report: K143467.txt
Output: k143467

Report: K171742.txt
Output: K171742

Report: K170491.txt
Output: K170491

Report: K171641.txt
Output: K171641

Report: K162526.txt
Output: k162526
```

**Wikipedia Setting** Next we provide the same demonstrations for the "born" attribute in the NBA Players Wikipedia dataset.

We provide examples of generated functions below:

```
def get_born_field(text: str):
    """
    Function to extract born.
    """
    pattern = r"Born(.*?)Nationality"
    result = re.findall(pattern, text, re.DOTALL)
    return result

def get_born_field(text: str):
    """
    Function to extract the born field.
    """
    soup = BeautifulSoup(text, parser="html.
     parser")
    born_field = soup.find('span', class_="bday")
    born_field = born_field.text
    return born_field

def get_born_field(text: str):
    """
    Function to extract born.
    """
    pattern = r"born in (.*?)\."
    matches = re.findall(pattern, text)
    return matches
```

We provide relevant snippets of three unique player documents from the dataset as examples for the case study.

```
Document: Luis_Flores.html
Content snippet: <th class="infobox-label">Born</
    th><td class="infobox-data"><span style="
    display:none"> (<span class="bday
    ">1981-04-11</span>) </span>April 11, 1981<
    span class="noprint ForceAgeToShow"> (age
    42)</span><br/>

Document: Magic_Johnson.html
Content snippet: <th class="infobox-label" scope
    ="row">Born</th><td class="infobox-data"><
    span style="display:none"> (<span class="
    bday">1959-08-14</span>) </span>August 14,
    1959<span class="noprint ForceAgeToShow"> (
    age 63)</span><br/>

Document: Draymond_Green.html
Content snippet: <th class="infobox-label" scope
    ="row">Born</th><td class="infobox-data"><
    span style="display:none"> (<span class="
    bday">1990-03-04</span>) </span>March 4,
    1990<span class="noprint ForceAgeToShow"> (
    age 32)</span><br/>
```

The resulting outputs of the functions applied to the three unique player documents from the dataset are shown below.

```
Document: Luis_Flores.html
Predictions: ['April 11, 1981', 'April 11, 1981',
     'April 11, 1981', 'April 11, 1981', 'April
    11, 1981', 'April 11', 'April 11',
    '1981-04-11', '(1981-04-11) April 11, 1981 (
    age\xa041)San Pedro de Macoris, Dominican
    Republic', '(1981-04-11) April 11, 1981 (age
    \xa041)San Pedro de Macoris, Dominican
    Republic']

Document: Magic_Johnson.html
Predictions: ['August 14, 1959', 'August 14,
    1959', 'August 14, 1959', '1959', '1959', '
    August 14', '1959, August 14, in <a href="/
    wiki/Lansing, child did not have HIV, 1956,
    to Melissa Mitchell. Although Andre was
    raised by his mother', '1959-08-14',
    '(1959-08-14) August 14, 1959 (age\xa063)
    Lansing, Michigan, U.S.', '(1959-08-14)
    August 14, 1959 (age\xa063)Lansing, Michigan
    , U.S.']

Document: Draymond_Green.html
Predictions: ['March 4, 1990', 'March 4, 1990', '
    March 4, 1990', '1990', '1990', 'March 4',
    '1990, March 4, with his then-girlfriend
    Jelissa Hardy.<sup class="reference" id="
    cite_ref-164"><a href="#cite_note
    -164">[164]</a></sup><sup class="reference"
    id="cite_ref-165"><a href="#cite_note
    -165">[165]</a></sup> In 2018, 2020.<sup
    class="reference" id="cite_ref-166"><a href
    ="#cite_note-166">[166]</a></sup><sup class
    ="reference" id="cite_ref-167"><a href="#
    cite_note-167">[167]</a></sup><sup class="
    reference" id="cite_ref-168"><a href="#
    cite_note-168">[168]</a></sup> They held
    their wedding ceremony on August 14',
    '1990-03-04', '(1990-03-04) March 4, 1990 (
    age\xa032)Saginaw, Michigan, U.S.',
    '(1990-03-04) March 4, 1990 (age\xa032)
    Saginaw, Michigan, U.S.']
```

For the examples above, the input vectors provided to the weak supervision model are as follows. Recall results the 10 outputs per document (obtained from the 10 candidate functions) are condensed to the top-$k$ (e.g. $k = 5$) most frequently occurring function outputs.

```
Document: Luis_Flores.html
Input: ['April 11, 1981', 'April 11',
    '1981-04-11', '(1981-04-11) April 11, 1981 (
    age\xa041)San Pedro de Macoris, Dominican
    Republic']

Document: Magic_Johnson.html
Inputs: ['August 14, 1959', '1959', August 14, in
     <a href="/wiki/Lansing, child did not have
    HIV, 1956, to Melissa Mitchell. Although
    Andre was raised by his mother',
    '1959-08-14', '(1959-08-14) August 14, 1959
    (age\xa063)Lansing, Michigan, U.S.']

Document: Draymond_Green.html
Inputs: ['March 4, 1990', '1990', 'with his then-
    girlfriend Jelissa Hardy.<sup class="
    reference" id="cite_ref-164"><a href="#
    cite_note-164">[164]</a></sup><sup class="
    reference" id="cite_ref-165"><a href="#
    cite_note-165">[165]</a></sup> In 2018,
    2020.<sup class="reference" id="cite_ref
    -166"><a href="#cite_note-166">[166]</a></
    sup><sup class="reference" id="cite_ref
    -167"><a href="#cite_note-167">[167]</a></
    sup><sup class="reference" id="cite_ref
    -168"><a href="#cite_note-168">[168]</a></
    sup> They held their wedding ceremony on
    August 14', '1990-03-04', '(1990-03-04)
    March 4, 1990 (age\xa032)Saginaw, Michigan,
    U.S.', '(1990-03-04) March 4, 1990 (age\
    xa032)Saginaw, Michigan, U.S.']
```

Applying the learned label model to the function predictions, the outputs are returned to the user by EVAPORATE are shown below. Note that this is an example of a difficult attribute with lower quality results.

```
Document: Luis_Flores.html
Output: 'april 11 1981'

Document: Magic_Johnson.html
Output: '1959 08 14 august 14 1959 age\
    xa063lansing michigan u.s.'

Document: Draymond_Green.html
Output: '1990 03 04 march 4 1990 age\xa032saginaw
     michigan u.s.'
```

Note that the cleaning steps discussed in Section F.2, could help convert these noisy outputs.

## G.2  Generated Functions

Here we provide additional examples of the candidate functions generated by EVAPORATE. Note that the displayed functions are randomly selected and not filtered for quality (some may fail or achieve low quality) to give a representative sampling.

First, within the FDA setting, the following functions are generated for the "intended use" attribute.

```python
def get_intended_use_field(text: str):
    """
    Function to extract intended use.
    """
    intended_use_field = re.findall(r'Intended
     Use(.*?)\n\n', text, re.DOTALL)
    return intended_use_field

def get_intended_use_field(text: str):
    """
    Function to extract intended use.
    """
    intended_use_field = re.findall(r'Intended
     Use\s*(.*?)\s*Similar', text, re.DOTALL)
    return intended_use_field

def get_intended_use_field(text: str):
    """
    Function to extract the intended use field.
    """
    parts = text.split("G. Intended Use:")[-1]
    intended_use_field = parts.split("F.
     Regulatory Information:")[0]
    return intended_use_field
```

Second, within the Wikipedia NBA setting, the following functions are generated for the "college" attribute.

```python
def get_college_field(text: str):
    """
    Function to extract the college field.
    """
    soup = BeautifulSoup(text, parser="html.
     parser")
    college_field = soup.find('th', string="
     College").find_next_sibling('td').text
    return college_field

def get_college_field(text: str):
    """
    Function to extract college.
    """
    pattern = r'<i>(.*?)</i>'
    college_field = re.findall(pattern, text)
    return college_field

def get_college_field(text: str):
    """
    Function to extract the college field.
    """
    pattern = r'College Basketball at Sports-
     Reference.com'
    college_field = re.findall(pattern, text)
    return college_field[0]

def get_college_field(text: str):
    """
    Function to extract college.
    """
    pattern = r"College.*?>(.*?)<"
    college_field = re.findall(pattern, text)
    return college_field
```

## H  SYSTEM INPUT AND OUTPUT DIAGRAMS

In Figures 7, 8, and 9, we include sample inputs and outputs for EVAPORATE.

## System Inputs



## System Output

| name | college name | position | height (ft) |
|---|---|---|---|
| Kevin Durant | Texas | Small forward / Power forward | 6 ft 10 in |
| Kermit Washington | American | Power forward | 6 ft 8 in |
| Rik Smiths | Marist | Center | 7 ft 4 in |
| Lamar Stevens | Penn State | Small forward | 6 ft 6 in |
| Tony Delk | Kentucky | Point guard / Shooting guard | 6 ft 2 in |
| Greg Foster | UTEP | Power forward / Center | 6 ft 11 in |
| Magic Johnson | Michigan State | Point guard | 6 ft 9 in |
| Tim Duncan | Wake Forest | Power forward / Center | 6 ft 11 in |

**Figure 7: Diagram depicting EVAPORATE input and sample output on the Wikipedia NBA Players (HTML) setting.**

## System Inputs



## System Output

| purpose for submission | 510(k) number | measurand | applicant | type of test |
|---|---|---|---|---|
| New assay | k141689 | C-reactive protein (CRP) | Qualigen, Inc. | Quantitative |
| New assay | k143502 | Opiates | Immunalysis Corporation | homogeneous enzyme immunoassay |
| New Device | k143075 | Sex Hormone Binding Globulin | Tosoh Bioscience | Quantitative immunoassay |
| New Device | k150168 | Tacrolimus | Siemens Healthcare | Quantitative immunoassay |
| New Device | k161714 | Barbiturates | Immunalysis Corporation | Homogenous Enzyme Immunoassay |
| New Device | k180209 | 1,5-Anhydroglucitol | Diazyme Laboratories Inc. | Colorometric, pyranose oxidase |
| New WSI System | K190332 | Not applicable | Leica Biosystems Imaging, Inc. | Digital pathology WS imaging |

**Figure 8: Diagram depicting EVAPORATE input and sample output on the Medical AI Device FDA Reports (TXT) setting.**

| System Inputs | System Output |
| --- | --- |

**System Inputs**

Watch HD Trailers on IMDb

GLOBES    List of Nominees   Photo Gallery   More Awards          Learn more »
More at IMDbPro »

**Sin City (2005)**

124 min  -  Crime | Drama | Thriller  -  1 April 2005 (USA)
1 2 3 4 5 6 7 8 9 10 8.3/10 X
Users: (264,723 votes) 1,651 reviews | Critics: 364 reviews

A film that explores the dark and miserable town, Basin City, and tells the story of three different people, all caught up in violent corruption.

Sin City Poster**Directors:**

Frank Miller, Robert Rodriguez, and 1 more credit »

**Writer:**

Frank Miller (graphic novels)

**Stars:**

Mickey Rourke, Clive Owen and Bruce Willis
Watch Trailer »

208 photos | 43 videos | 2220 news articles »
**Top 250 #101** | 16 wins & 29 nominations See more awards »
Edit

**Cast**

Cast overview, first billed only:

Jessica Alba          ... Nancy Callahan

Devon Aoki          ... Miho

**System Output**

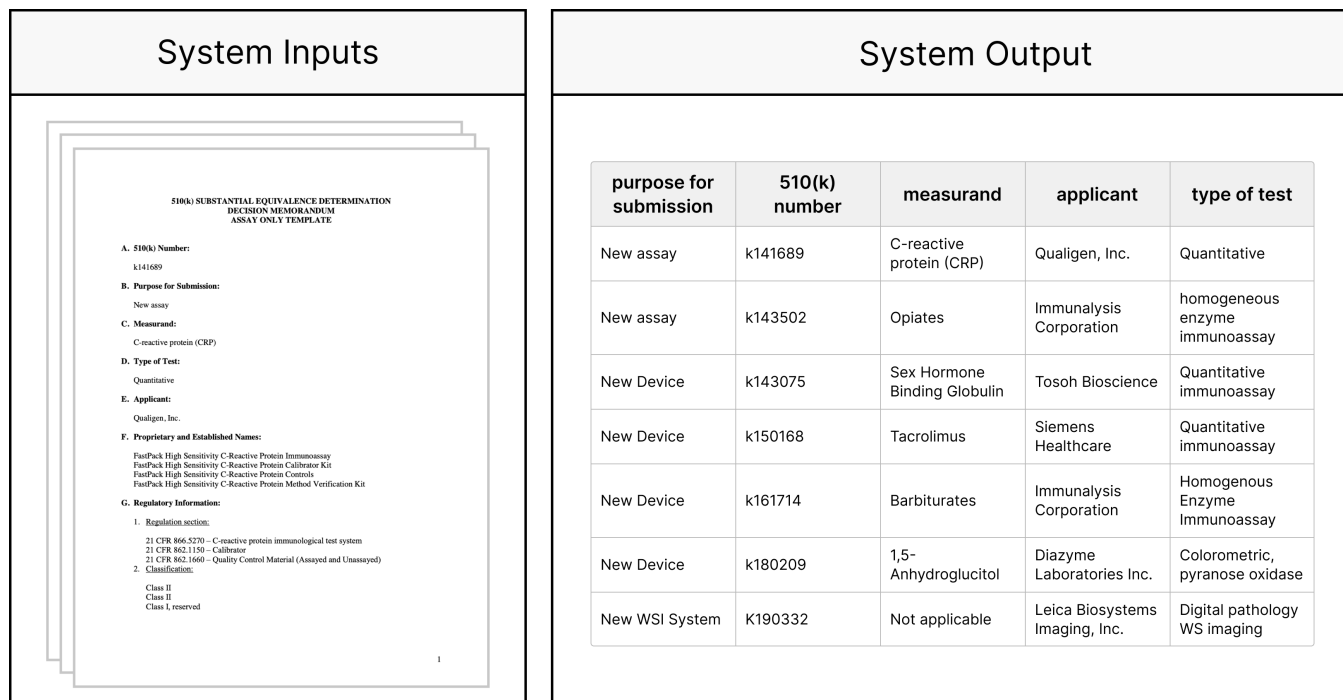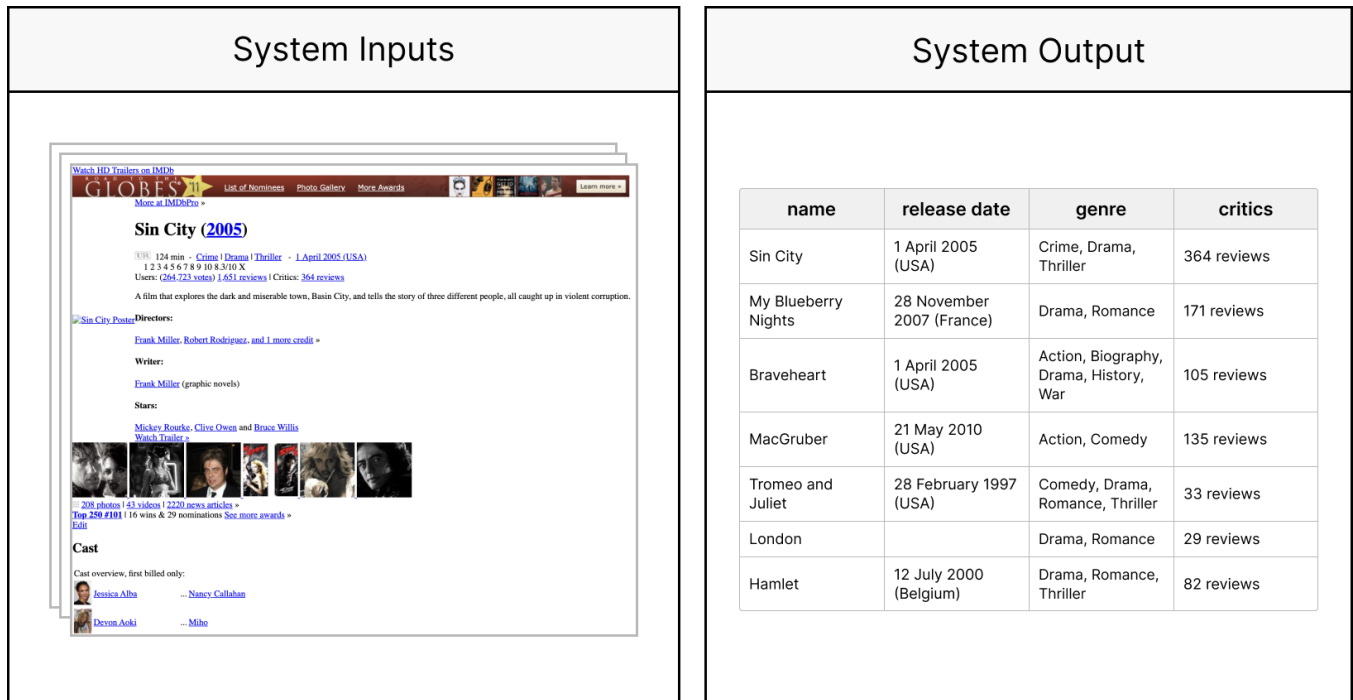| name | release date | genre | critics |
| --- | --- | --- | --- |
| Sin City | 1 April 2005 (USA) | Crime, Drama, Thriller | 364 reviews |
| My Blueberry Nights | 28 November 2007 (France) | Drama, Romance | 171 reviews |
| Braveheart | 1 April 2005 (USA) | Action, Biography, Drama, History, War | 105 reviews |
| MacGruber | 21 May 2010 (USA) | Action, Comedy | 135 reviews |
| Tromeo and Juliet | 28 February 1997 (USA) | Comedy, Drama, Romance, Thriller | 33 reviews |
| London | | Drama, Romance | 29 reviews |
| Hamlet | 12 July 2000 (Belgium) | Drama, Romance, Thriller | 82 reviews |

**Figure 9: Diagram depicting EVAPORATE input and sample output on the SWDE Movies IMDB (HTML) setting.**