# Multi-Agent AI Framework for Automated Corporate Data Enrichment

**SAHIN UREGIL[1*], YAGIZ ERDEM[1*], SULTAN TURHAN[2], GÜNCE K. ORMAN[2], and DOGA Y. YILMAZ[1],**

[1]Technology & Innovation, Kariyer.net, Istanbul 34768, Turkey (e-mail: sahin.uregil@kariyer.net, yagiz.erdem@kariyer.net, doga.yilmaz@kariyer.net)
[2]Department of Computer Engineering, Galatasaray University, Istanbul 34349, Turkey (e-mail: sturhan@gsu.edu.tr, korman@gsu.edu.tr)
[*]These authors contributed equally to this work.

Corresponding authors: Sahin Uregil (e-mail: sahin.uregil@kariyer.net) and Yagiz Erdem (yagiz.erdem@kariyer.net).

**ABSTRACT** The ability to rapidly acquire and process comprehensive corporate information has become critical for market analysis and competitive intelligence. Manual extraction and structuring of business data from web sources remains time-intensive and resource-demanding, limiting business intelligence scalability. We developed an automated company data enrichment system that transforms basic company identifiers into structured business intelligence profiles using multi-agent artificial intelligence (AI) frameworks and large language models (LLMs). The system implements a four-stage pipeline comprising query generation, automated information gathering, structured data synthesis, and quality validation mechanisms. The approach was evaluated on a dataset of 400 companies sourced from Kariyer.net's Turkish and global company listings. The system demonstrated varying extraction success rates across information categories: founding year (86.8%), parent company information (64.0%), subsidiary information (52.8%), funding summaries (52.0%), and founder names (22.8%). Analysis revealed significant organizational complexity within the Turkish corporate landscape, with conglomerate companies averaging 25.44 subsidiaries each compared to 6.45 for all companies. The system provides a scalable solution for enterprise-level business intelligence applications, enabling organizations to transform basic company identifiers into comprehensive corporate profiles with minimal manual intervention. The framework reduces market research timelines from weeks to hours while revealing important patterns in organizational complexity within emerging markets.

**INDEX TERMS** Artificial intelligence, automated data processing, business intelligence, corporate information extraction, data enrichment, hierarchical databases, knowledge extraction, large language models, multi-agent systems, web scraping

## I. INTRODUCTION

In the contemporary business landscape, the ability to rapidly acquire and process comprehensive corporate information has become a critical competitive advantage for organizations engaged in market analysis, competitive intelligence, and strategic decision-making (1; 2). The proliferation of digital information sources has created vast repositories of unstructured business data distributed across corporate websites, news articles, and public databases (1). However, the manual extraction and structuring of this information remains a time-intensive and resource-demanding process that significantly limits the scale and efficiency of business intelligence operations (1; 2; 3).

The primary objective of this study is to develop and validate an automated company data enrichment system that transforms basic company identifiers into comprehensive, structured business intelligence profiles. This system leverages multi-agent AI frameworks (4; 5; 6) and LLMs (7; 8; 9) to systematically extract, process, and structure corporate information from web-based sources. The proposed solution aims to bridge the gap between raw company naming data and the rich, multi-dimensional corporate profiles required for advanced business analytics, competitive analysis, and market intelligence applications.

Several critical challenges emerge when attempting to automate corporate information extraction at scale. First, the heterogeneous nature of web-based corporate information presents significant data quality and consistency issues, as companies vary widely in their information disclosure practices and digital presence (1; 10; 11; 12). Second, complex

organizational structures, particularly within conglomerate companies, require sophisticated parsing and relationship mapping capabilities to accurately represent subsidiary networks and parent-child relationships (13). Third, the semantic variability in how corporate information is presented across different sources necessitates robust natural language processing (NLP) approaches capable of identifying and extracting relevant information regardless of presentation format (14). Fourth, the need for structured output that conforms to predefined schemas while maintaining flexibility for diverse business intelligence applications presents additional technical complexity (15).

Our study addresses these challenges through the development of a comprehensive multi-stage pipeline that combines query generation, automated web information gathering, structured data synthesis, and quality validation mechanisms. By implementing this automated approach on a dataset of 400 companies sourced from Kariyer.net's[1] database of Turkish and global company listings, we demonstrate the system's capability to extract diverse corporate attributes with varying success rates, revealing important patterns in information availability and organizational complexity within the Turkish business landscape. The resulting framework provides a scalable solution for enterprise-level business intelligence applications while offering insights into the challenges and opportunities associated with automated corporate data enrichment in emerging markets.

The primary contributions of this study to the field of automated business intelligence and corporate data analysis are as follows:

- **Comprehensive Corporate Profiling Framework:** We develop an end-to-end system that transforms basic company identifiers into detailed and meaningful multi-dimensional business profiles using LLMs, enabling organizations to rapidly assess market opportunities without extensive manual research.
- **Automated Hierarchical Relationship Discovery:** Our system automatically identifies and structures complex parent-subsidiary relationships within corporate ecosystems.
- **Up-to-date Market Intelligence through Web Automation:** By leveraging automated web data collection (16), our system provides access to current market information, enabling businesses to make data-driven decisions based on up-to-date intelligence.
- **Scalable Business Intelligence Infrastructure:** The system generates machine-readable JavaScript Object Notation (JSON)[2] outputs that integrate seamlessly with existing business intelligence platforms and analytical workflows, eliminating manual data entry bottlenecks and enabling automated downstream processing for large-scale market analysis.

---

[1]https://www.kariyer.net
[2]https://www.json.org/json-en.html

- **Quality-Assured Information Extraction:** We implement a reflection mechanism that automatically validates extracted information quality and triggers reprocessing when necessary, ensuring reliable business intelligence outputs.

The rest of the paper is organized as follows. Section II discusses the materials and methods used in the development of the automated corporate data enrichment system, including the dataset, query generation, information gathering, and data structuring mechanisms. Section III presents the experimental setup and results, covering the system architecture, parameter configuration, and performance analysis. It also includes the results of the data extraction process, with detailed insights into feature completion rates and organizational complexity, a representative running example, and an evaluation of the business impact and operational value generated by the system. Section IV concludes the paper by discussing related work and suggesting potential directions for future research and system enhancements.

## II. MATERIALS AND METHODS
### A. DATASET OVERVIEW

The dataset utilized in this study comprises a curated list of company names primarily sourced from Kariyer.net, a leading employment and recruitment platform in Turkey. The selection includes a broad range of companies operating in various sectors, with an emphasis on both Turkish enterprises and multinational corporations that have a significant presence in the Turkish market.

The initial dataset consisted of 400 company names sampled directly from Kariyer.net's internal listings. These include a diverse mix of small, medium, and large-scale organizations, with a subset of the entries (45 companies) identified as conglomerates. The inclusion of conglomerates provides an added dimension of structural complexity due to their typically multi-sector and multi-brand configurations (17).

The dataset was constructed in close collaboration with domain experts from Kariyer.net, ensuring that the selected companies reflect a representative sample of the broader employment landscape within Turkey. The primary data field at this stage is the *company name*, which serves as the foundational identifier for all subsequent processing and analysis.

No additional attributes or metadata (e.g., business activities, relationships, or organizational structures) were found in the raw dataset. As such, the dataset initially served as a naming index, intended to be further enhanced with detailed, meaningful and structured information in later stages of the project.

### B. PROBLEM STATEMENT

While the initial dataset provided a solid starting point through its comprehensive list of company names, it lacked the depth and contextual richness required for meaningful business analysis. Company names alone do not offer sufficient semantic information about a company's operations,

industry focus, or organizational structure (18). This limitation makes it difficult to derive actionable insights or perform tasks such as clustering, segmentation, or relationship mapping with any degree of accuracy or business relevance (19).

One of the primary challenges stemmed from the absence of descriptive metadata that could convey the nature of each company's products, services, or sectoral alignment. Without such information, any attempt to group or compare companies based on their business activities would be inherently limited and potentially misleading (19).

Another significant challenge involved conglomerate companies, which often operate across multiple industries through various subsidiaries or sub-brands. Treating these entities as monolithic units can introduce substantial noise into any analytical process, obscuring important intra-organizational distinctions and misrepresenting inter-company relationships. Disaggregating these complex structures requires an understanding of subsidiary relationships that was not present in the original dataset (20).

Furthermore, the lack of auxiliary information, such as founding year, founder names, or hierarchical relationships, restricted the potential for cross-company linkages or historical and structural analysis. These gaps not only hindered the development of a nuanced understanding of the corporate landscape but also limited the dataset's applicability for advanced business intelligence tasks (19).

In sum, the initial dataset's narrow focus on company names posed substantial barriers to achieving the semantic and structural granularity needed for high-quality, business-driven analysis.

To mitigate these problems, a multi-agent AI system has been devised to systematically enrich and structure the dataset. This system performs detailed and meaningful enrichment of corporate metadata by automatically extracting, contextualizing, and structuring diverse business attributes from web-based sources, transforming sparse company identifiers into comprehensive organizational profiles. A flowchart outlining the overall architecture of the proposed system is presented in Fig. 1.

### C. QUERY GENERATION

The query generation phase constitutes the initial step of the automated information extraction pipeline. This component takes two primary inputs: the *company name*, which serves as the subject of investigation, and an *extraction schema*, which specifies the structured format and target attributes for data retrieval.

The *extraction schema* functions as a flexible blueprint guiding the downstream stages of the system. It defines the detailed structure of the desired output, allowing users to tailor the pipeline to varying business or analytical needs. The schema is represented in a JSON-like format, where each field corresponds to a particular dimension of company-related information, such as founding year, product descriptions, founders, or organizational affiliations. This standardized representation facilitates seamless data exchange and

interoperability with existing enterprise systems, enabling straightforward integration into diverse business intelligence ecosystems without requiring custom data transformation processes. The modularity of the schema enables the system to dynamically adjust its behavior based on user-specified informational requirements.

For the purposes of this study, a default schema to support general business analysis tailored to our internal business use case was developed. The schema encompasses seven primary fields designed to capture essential business intelligence attributes: `company_name` (official company designation), `founding_year` (year of establishment), `founder_names` (founding team members), `product_description` (main goods or services offered), `funding_summary` (investment and funding history), `branch_dealer` (regional branches or dealer networks), and `affiliations` (organizational relationships including `subsidiaries` and `parents`).

Given these two inputs, *company name* and *extraction schema*, the system proceeds to generate a fixed number of web search queries using an LLM according to *max search queries* parameter as detailed in `QUERY_GENERATOR` row of the Table 1. These queries are semantically aligned with the schema fields and are designed to maximize the likelihood of retrieving relevant and structured information from publicly available sources.

The queries represent the information-seeking intent derived from the schema and set the stage for the subsequent phase of information gathering, where the system uses these queries to interact with web search engines and collect relevant documents.

### D. INFORMATION GATHERING

Following the generation of targeted search queries, the next phase in the pipeline involves automated information gathering from the web. This step aims to collect and prepare high-quality textual data that corresponds to each of the previously formulated queries.

The process begins by feeding the generated queries into a web automation module that interacts with a publicly available search engine. For every individual query, the system retrieves a total of $2 \times$ `max_search_results` URLs. The rationale behind collecting twice the target number of URLs is to maintain a buffer in case some of the pages are unreachable.

After URL collection, the HyperText Markup Language (HTML)[3] content of each page is downloaded. These raw web documents are then parsed to extract the primary text content, stripping away non-informative elements such as scripts, styles, and advertisements.

Each cleaned textual document is paired with its corresponding query and passed to an LLM. The LLM is tasked with identifying and extracting only the content that directly answers the given query, effectively acting as a semantic filter. This ensures that only relevant, concise, and contextually

---

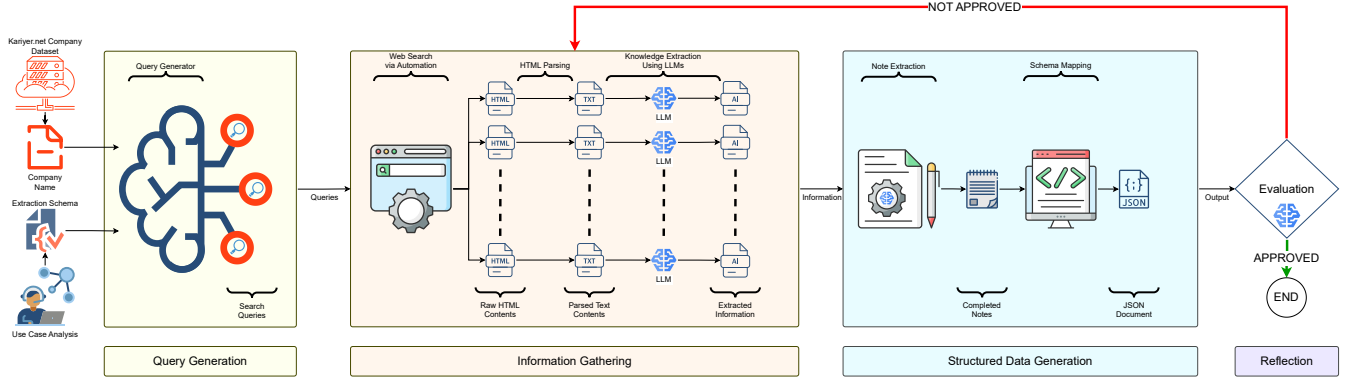[3]https://html.spec.whatwg.org/

FIGURE 1: System Architecture for Automated Company Data Enrichment and Structuring

TABLE 1: Overview of System Prompts for Company Information Extraction Pipeline

| Prompt Name | Purpose | Input Variables | Output Format |
|---|---|---|---|
| QUERY_GENERATOR | Generate targeted search queries to gather factual company data | `{company_name}`,`{max_search_queries}`,`{extraction_schema}` | List of strings |
| PAGE_SUMMARIZER | Clean and summarize page content in response to a specific query | `{page_content}`,`{query}` | Unstructured text |
| NOTE_EXTRACTION | Extract and organize notes about a company based on raw web content | `{company_name}`,`{extraction_schema}`,`{content}` | Markdown file |
| SCHEMA_MAPPING | Convert unstructured research notes into structured data as per schema | `{extraction_schema}`,`{notes}` | Structured text (JSON-like) |
| REFLECTION | Evaluate completeness and quality of the structured information against the schema | `{schema}`,`{extraction_schema}` | Boolean value |

grounded information is carried forward to the next stage of the pipeline.

At the conclusion of this phase, the system produces up to `max_search_queries × max_search_results` units of cleaned and semantically filtered textual content. These documents form the raw material for the structured data generation process that follows. The input and output structure of this step is given in detail in `PAGE_SUMMARIZER` row of the Table 1.

### E. STRUCTURED DATA GENERATION
Once query-specific content has been extracted and filtered, the next stage in the pipeline involves synthesizing this information into a structured and standardized format. This step bridges the gap between unstructured textual data and the structured outputs required for downstream analysis.

The process begins with a summarization task, in which an LLM processes all cleaned content collected for a single company across multiple queries. The LLM consolidates relevant information into a unified markdown-style research note, guided by the predefined *extraction schema*. This summary captures all schema-relevant insights discovered during the

earlier information gathering phase, including partial or ambiguous information as explained in `NOTE_EXTRACTION` row of the Table 1. The summarization agent also flags areas of uncertainty or inconsistency across sources, supporting transparency and further human validation where necessary.

Following summarization, the research note is passed to the schema mapping agent, which converts the narrative content into a machine-readable structured format that complies with the original *extraction schema*. This final structured output is expressed in a JSON-style format as detailed in `SCHEMA_MAPPING` row of the Table 1, enabling integration with business intelligence tools or analytical pipelines.

This structured output not only enables consistency across companies but also supports automated processing, querying, and large-scale analysis. The integration of LLMs at this stage significantly reduces manual effort while maintaining interpretability and adaptability to evolving information needs.

### F. REFLECTION
The final step in the pipeline involves a validation mechanism designed to assess the reliability and completeness of the structured output. This step, referred to as *Reflection*, employs

an LLM to inspect the generated JSON document against the original *extraction schema* and *completed notes* and determine whether the information appears consistent, comprehensive, and plausible as the LLM employs meta reasoning to act as a semantic filter, identifying and extracting only the content that directly answers the given *company name* and *extraction schema* while evaluating its relevance and completeness (21).

The LLM produces a boolean flag called `is_satisfactory`. If this flag is set to `True`, the process terminates, indicating that the output is deemed acceptable. If it is set to `False`, and the parameter `max_reflection_steps` is greater than zero, the system re-initiates the information gathering phase and attempts to improve the result. This mechanism acts as a safeguard against hallucinations and incomplete outputs, improving the trustworthiness of the autonomous pipeline in a similar fashion to *chain of thought* (22; 23; 24; 25).

## III. EXPERIMENTS AND RESULTS

### A. EXPERIMENTAL SETUP

The experimental configuration was designed to balance comprehensiveness with computational efficiency while ensuring high-quality data extraction. This section details the key parameters, schema definition, and technical specifications employed throughout the experimental process.

The automated company data enrichment system developed for this study builds upon the foundational architecture of the *company-researcher* framework from the LangChain[4] repository[5]. While this open-source framework provided the initial structural blueprint for web-based company information extraction, substantial modifications and enhancements were implemented to address the specific requirements of our research objectives and the Turkish business landscape.

### 1) Extraction Schema Configuration

The core extraction schema captures essential business intelligence attributes and encompasses seven primary fields as detailed in II-C. The `founding_year` and `funding_summary` fields are designed to be nullable (i.e., allowed to be `null` if information is unavailable). Conversely, the `founder_names`, `branch_dealer`, `parents` (within `affiliations`), and `subsidiaries` (within `affiliations`) fields are designed to allow empty lists (`[]`) when no corresponding information is found, rather than being strictly null. This distinction allows for explicit representation of the absence of data.

Furthermore, this extraction schema is formally defined using Pydantic[6] classes within the system's codebase, which provides robust runtime validation and ensures strict adherence to the specified data types and nullability/emptiness rules for each field, mirroring the JSON schema provided to the LLM.

The complete schema definition follows a JSON schema format with explicit type constraints and validation rules as illustrated in Fig. 2.
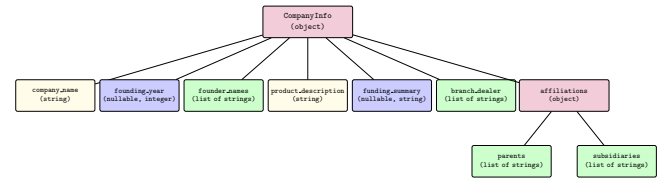


FIGURE 2: Company Information Extraction Schema Structure

### 2) System Parameters and Configuration

The experimental setup employed a calibrated set of parameters to optimize data quality and processing efficiency. The query generation phase was configured to produce eight search queries per company (`max_search_queries=8`), providing comprehensive schema coverage while maintaining manageable computational overhead.

For each query, the system retrieved up to three search results (`max_search_results=3`), resulting in a maximum of 24 candidate web pages per company. A character limit of 250,000 characters per web page was enforced (`max_character_count_for_one_page=250000`), equivalent to approximately 62,500 tokens, ensuring stable LLM performance.

The system implements parallel processing through batch processing with a batch size of four (`batch_size=4`), enabling concurrent query processing in two sequential batches. The reflection mechanism was configured with a single iteration limit (`max_reflection_steps=1`) to balance quality assurance with computational efficiency.

### 3) Technical Implementation

The web automation component utilizes Selenium WebDriver[7] for browser-based interactions as it is among the most popular web automation tools (26), enabling robust handling of dynamic web content. The HTML parsing pipeline combines LXML[8] and BeautifulSoup[9] libraries for reliable content extraction across diverse web page structures. Throughout the entire pipeline, *Gemini-2.0-flash*[10] serves as the LLM for query generation, content filtering, summarization, schema mapping, and reflection as it was found to be adequately well-performing according to the existing research (27; 28; 29). The system architecture supports both single-company processing and bulk operations, with integrated error handling and logging throughout the extraction process.

---

[4]https://www.langchain.com/
[5]https://github.com/langchain-ai/company-researcher
[6]https://docs.pydantic.dev/latest/

[7]https://www.selenium.dev/
[8]https://lxml.de/
[9]https://www.crummy.com/software/BeautifulSoup/
[10]https://gemini.google.com/app

## B. DATA EXTRACTION PERFORMANCE ANALYSIS

The automated company data enrichment system demonstrated varying levels of success across different schema fields when applied to the 400-company dataset. This section presents a comprehensive analysis of feature completion rates and structural complexity patterns discovered through the extraction process.

### 1) Feature Completion Rates

The system's performance varied significantly across different information categories, as illustrated in Fig. 3. The most successfully extracted field was the founding year, with 347 out of 400 companies (86.8%) having this information successfully identified and structured. This high completion rate reflects the widespread availability and standardized reporting of establishment dates across corporate websites and public records.

Corporate relationship data showed moderate extraction success, with parent company information successfully identified for 256 companies (64.0%) and subsidiary information for 211 companies (52.8%). The relatively high success rate for parent company identification suggests that ownership structures are frequently disclosed in corporate communications.

Funding-related information presented moderate extraction challenges, with funding summaries successfully compiled for 208 companies (52.0%). This completion rate reflects the heterogeneous nature of funding disclosure practices, where privately held companies often maintain limited public information about their financial backing, while venture-backed and publicly funded entities typically provide more comprehensive funding histories.

The most challenging field proved to be founder names extraction, with only 91 companies (22.8%) having complete founder information successfully identified. This low completion rate can be attributed to several factors: historical founders may not be prominently featured, founding team information is often relegated to company history sections that may be difficult to locate, and many established companies focus their public communications on current leadership rather than founding members.

These statistics, illustrated in Fig. 3, are provided specifically for features that are either nullable (i.e., `founding_year`, `funding_summary`) or allowed to be empty lists (i.e., `founder_names`, `branch_dealer`, `parents`, `subsidiaries`). Other features within the schema, such as `company_name` and `product_description`, are always expected to be filled, and thus their completion rates are inherently 100%.

### 2) Subsidiary Structure Analysis

The analysis of subsidiary structures revealed significant complexity within the Turkish corporate landscape, as depicted in Fig. 4. Across all 400 companies, the system identified an average of 6.45 subsidiaries per company, with substantial variation in organizational complexity.
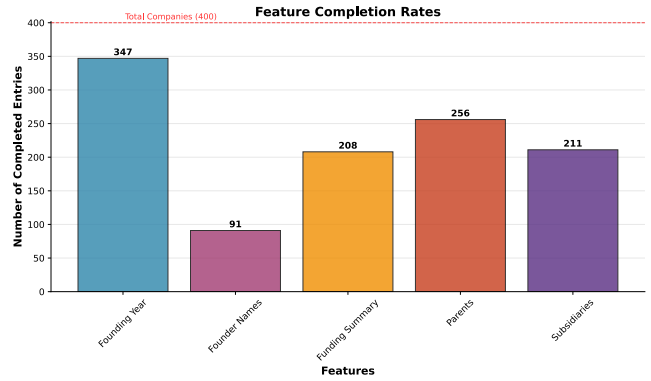


FIGURE 3: Feature Completion Rates

The subsidiary distribution exhibits a pronounced right-skewed pattern, with 189 companies (47.3%) having no identified subsidiaries and 211 companies (52.8%) maintaining at least one subsidiary relationship. The maximum subsidiary count reached 89 for a single organization, demonstrating the extreme complexity possible within conglomerate structures. The standard deviation of 13.34 reflects this high variability in organizational complexity across the dataset.

Conglomerate companies, representing 45 entities within the dataset, displayed markedly different structural characteristics compared to the general population. These organizations averaged 25.44 subsidiaries each, nearly four times the overall average, confirming their role as complex multi-entity corporate structures. This finding validates the initial hypothesis that conglomerates would present unique challenges for organizational analysis and require specialized handling in business intelligence applications.

The subsidiary analysis reveals a bimodal corporate landscape in Turkey, characterized by a large number of relatively simple organizational structures alongside a smaller subset of highly complex conglomerate entities. This distribution pattern has important implications for business analysis methodologies, suggesting that different analytical approaches may be required for different organizational complexity levels within the Turkish corporate environment.
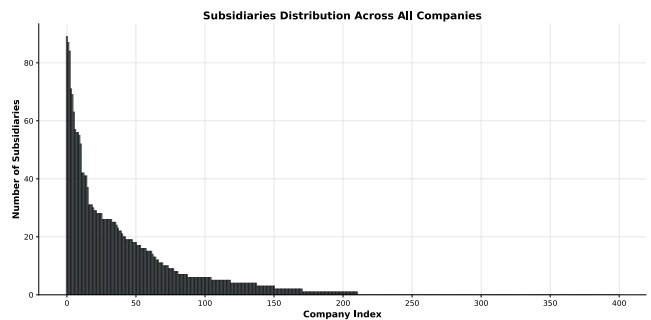


FIGURE 4: Subsidiaries Distribution

## C. RUNNING EXAMPLE

To better illustrate the workings of this process throughout the explained steps, we follow a concrete running example: *Kariyer.net A.Ş.*, a prominent employment and recruitment company in Turkey. The following subsections outline each phase of the pipeline with their respective inputs and outputs.

### Query Generation

*Inputs:* The system takes the following inputs:

- `company_name: Kariyer.net A.Ş.`
- `max_search_queries: 8`
- `extraction_schema: Schema with user defined fields`

*Outputs:* The system produces 8 semantically-aligned search queries:

- `Kariyer.net A.Ş. founding year`
- `Kariyer.net A.Ş. founder names`
- `Kariyer.net A.Ş. product description`
- `Kariyer.net A.Ş. subsidiaries`
- `Kariyer.net A.Ş. parent company`
- `Kariyer.net A.Ş. company profile`
- `Kariyer.net A.Ş. branch dealers`
- `Kariyer.net A.Ş. about`

### Information Gathering

*Inputs:* For each of the 8 queries generated previously, the system performs a web search with:

- `query: (e.g., Kariyer.net A.Ş. about)`
- `max_search_results: 3`

*Outputs:* Up to 3 unstructured textual contents per query are collected via web automation and filtered by an LLM, resulting in a maximum of 24 cleaned page contents. Below is one such example for the query *"Kariyer.net A.Ş. about"*:

```
Kariyer.net is Turkey's largest employment
platform, operating since 1999. It aims to
connect job seekers and employers,
offering job postings and career guidance.
It also provides resources for university
selection and career planning. Kariyer.net
offers various filters for job searching,
enabling users to find opportunities based
on city, sector, position level,
department, work style, and education
level. The platform provides guidance at
every stage of the career journey and
provides career advice on resume
preparation, job searching, interviews,
and work life. Kariyer.net also offers
career resources such as the Position
Guide, University Guide, and Career Guide.
The HTML content also includes information
about Kariyer.net being a part of a group
of companies. The listed companies are:
arabam.com, chemOrbis, cimri.com,
emlakjet, hangikredi.com, neredekal.com,
sigortam.net, and steelOrbis.
```

### Structured Data Generation

*Note Extraction Phase*

*Inputs:*

- `company_name: Kariyer.net A.Ş.`
- `extraction_schema: Extraction schema used in the previous step`
- `content: Combined content of the 24 cleaned documents from previous step`

*Outputs:* A markdown-style research note is generated in Fig. 5:

*Schema Mapping Phase*

*Inputs:*

- `extraction_schema: Extraction schema used in the previous step`
- `notes: Research note generated in previous phase`

*Outputs:* A structured JSON-style company profile is created in Fig. 6:

```json
{
  "company_name": "Kariyer.net A.Ş.",
  "founding_year": 1999,
  "founder_names": [],
  "product_description": "Turkey's largest
  →  online employment platform, connecting
  →  job seekers and employers. Offers job
  →  postings, career guidance, resume
  →  database access, and tools for both job
  →  seekers and employers across multiple
  →  industries. Provides resources for
  →  university selection and career
  →  planning.",
  "branch_dealer": [],
  "affiliations": {
    "parents": [
      "iLab Holding"
    ],
    "subsidiaries": [
      "İşin Olsun",
      "İşkolig",
      "Techcareer.net",
      "WeBrand"
    ]
  }
}
```

FIGURE 6: Structured JSON-style output generated by the system.

### Reflection

In the case of our running example, the returned `is_satisfactory` flag from the reflection step was *True*, indicating that the extracted information met the completeness and quality criteria defined by the schema. As a result, the pipeline halted further iterations and finalized the structured output without requiring re-processing.

To further ensure the reliability and correctness of the reflection mechanism, we conducted a thorough evaluation of the system's outputs in collaboration with domain experts from Kariyer.net. Specifically, we inspected both the structured outputs and their corresponding reflection results for both this running example and the entire dataset of 400 companies. This joint review process confirmed that the reflection step effectively identified incomplete or inconsistent information and successfully triggered re-processing when necessary. These observations validate the reflection mechanism's role in enhancing output quality and reinforce the overall system's trustworthiness for enterprise-grade data enrichment tasks.

## D. BUSINESS IMPACT AND OPERATIONAL VALUE CREATION

To evaluate the computational efficiency of our system, we conducted timing analyses using free-tier resources, as shown

```
# Kariyer.net A.Ş. Research Notes
- **Company Name:** Kariyer.net A.Ş.
- **Founding Year:** 1999
- **Founder Names:** Not found in provided sources.
- **Product Description:** Turkey's largest online employment platform, connecting job seekers and
↪   employers. Offers job postings, career guidance, resume database access, and tools for both job
↪   seekers and employers across multiple industries. Provides resources for university selection and
↪   career planning.
- **Branch Dealers:** Kariyer.net has offices in Istanbul Ümraniye and Yeşilköy, Adana, Ankara,
↪   Antalya, Bursa, Denizli, Elazığ, Eskişehir, Gaziantep, İzmir, Kayseri, Kocaeli, Konya, Mersin,
↪   Sakarya, and Samsun.
- **Affiliations:**
  - **Parents:** iLab Holding (established in 2000 by Mustafa Say through Access Turkey Capital
  ↪   Group); Kariyer.net became part of iLab in 2006 (acquired on March 8, 2006).
  - **Subsidiaries:**
    - İşin Olsun (Launched in 2017)
    - İşkolig (Added in 2018)
    - Techcareer.net (Launched in 2021)
    - WeBrand (Launched in 2021)
- **Investments Made By Kariyer.net:** PeopleBox (Later Stage VC, 29-Apr-2022), Coensio (Corporate,
↪   07-Sep-2021), İşkolig (Buyout/LBO, 27-Jun-2018)
- **Companies within the same group as Kariyer.net:** arabam.com, chemOrbis, cimri.com, emlakjet,
↪   hangikredi.com, neredekal.com, sigortam.net, steelOrbis
- **Investors:** iLab Ventures (Venture Capital investor holding a majority stake), AccessTurkey
↪   Capital Group
- **Competitors:** Company-A, Company-B, Company-C, Company-D, Company-E, Company-F, Company-G,
↪   Company-H, Company-I, Company-J, Company-K
- **Headquarters:** Saray Mahallesi Site Yolu Sokak Anel İş Merkezi No: 5 Kat: 3 Istanbul, 34768,
↪   Turkey.
- **R&D Center Registration Date:** January 1, 2014
- **Membership Date (ARGEMİP):** October 18, 2019
```

FIGURE 5: Structured markdown-style output generated by the system. Company names have been anonymized to avoid the use of trademarks in accordance with academic publishing standards. These placeholders represent actual competitors identified by the system and correspond to various employment and recruitment platforms.

in Table 2. The system ultimately achieved an average processing time of 169.78 seconds per company across the complete dataset of 400 companies, which would normally take weeks if performed manually. As illustrated in the table, the total processing time increases approximately linearly with the number of companies, indicating that the system maintains consistent throughput and scales predictably with dataset size.

TABLE 2: System Processing Time Analysis Across Different Dataset Sizes

| Dataset Size | Total Time (s) | Average Time (s) |
|---|---|---|
| 25 | 4,435.37 | 177.41 |
| 50 | 8,730.46 | 174.61 |
| 100 | 17,411.34 | 174.11 |
| 200 | 34,156.47 | 170.78 |
| 400 | 67,911.52 | 169.78 |

The automated company data enrichment system delivers substantial business value through its ability to transform basic company identifiers into comprehensive, structured business intelligence. This transformation addresses critical operational challenges faced by organizations operating in the Turkish business ecosystem, particularly in the context of market analysis and strategic decision-making.

The system's ability to systematically extract and struc-

ture company information enables organizations to develop comprehensive market intelligence with significantly reduced manual effort. Traditional market research approaches requiring weeks of manual data collection and analysis can now be accomplished in hours through automated processing.

The automated extraction of product descriptions and service offerings facilitates rapid competitive benchmarking and market positioning analysis. Organizations can now efficiently identify direct competitors, and assess competitive threats without extensive manual research.

The automated nature of the data enrichment process enables scalable business development activities that were previously constrained by manual research limitations. Sales and business development teams can now rapidly assess and prioritize potential clients based on comprehensive company profiles, founding characteristics, and organizational complexity metrics.

The system's ability to process large sets of companies efficiently, demonstrates its utility for large-scale market analysis. Organizations can now systematically evaluate entire market segments, identify high-potential prospects based on specific criteria (such as founding year, organizational complexity, or parent company relationships), and optimize resource allocation for business development activities.

## IV. DISCUSSION

The key contributions of this study, ranging from **Comprehensive Corporate Profiling Framework** to **Automated Hierarchical Relationship Discovery**, were not just theoretical aspirations but were actively demonstrated through practical experiments. Each component of the system was subjected to real-world evaluation using a 400-company dataset, ensuring that the proposed solutions translated into measurable outcomes. For instance, the generation of detailed, machine-readable company profiles from only company names (as demonstrated in Section III-C) validates our **Comprehensive Corporate Profiling Framework**. Similarly, the automated extraction of hierarchical relationships (Fig. 4) empirically demonstrates the system's ability to model complex corporate ecosystems, validating the **Automated Hierarchical Relationship Discovery** contribution.

Further experimental validation was provided through live web automation (*Information Gathering* phase of Section III-C), where the system consistently acquired current and relevant data across diverse sources, demonstrating its support for **Up-to-date Market Intelligence through Web Automation**. The outputs were automatically structured into standardized JSON format (*Structured Data Generation* phase of Section III-C), reinforcing the system's **Scalable Business Intelligence Infrastructure**. Lastly, the reflection mechanism (*Reflection* phase of Section III-C) enabled quality assurance by re-triggering incomplete or inconsistent pipelines, which confirms the system's **Quality-Assured Information Extraction** capability. These experiments collectively anchor our contributions in applied results, reinforcing the system's readiness for business-oriented applications.

The automated company data enrichment system presented in this study offers a novel approach to transforming unstructured web-based information into actionable business intelligence for the Turkish corporate environment. Our findings demonstrate the system's efficacy in extracting diverse attributes, including company name, founding year, product descriptions, and affiliations, and highlight the varying success rates depending on the attribute's prominence and consistency across public web sources. The analysis of subsidiary structures further underscores the system's capability to untangle complex corporate hierarchies, particularly within conglomerates.

This work aligns with recent advancements in leveraging LLMs and agent-based frameworks for enterprise information processing. Similar to EICopilot (30), a solution enhancing search and exploration over large-scale knowledge graphs, our system employs LLM-driven agents to automate information extraction and summarization from natural language queries. Both systems recognize the inherent complexity of enterprise data and strive to improve efficiency and accuracy compared to manual methods. Furthermore, the incorporation of a "reflection" mechanism in our pipeline, designed to assess the reliability and completeness of the structured output and trigger re-iterations if unsatisfactory, mirrors the iterative refinement processes seen in EICopilot

and the annual report analysis method, highlighting a shared best practice in building robust LLM-driven data extraction systems.

However, a key distinction lies in the primary focus. While EICopilot is engineered to navigate and query existing, pre-structured knowledge graphs, leveraging technologies like Gremlin scripts[11] for precise data retrieval, our system adopts a more foundational and streamlined methodology that addresses the challenge of initial data enrichment and structuring through transforming raw, unstructured web content (including natural language descriptions and HTML content) into organized, machine-processable company profiles. This 'structuring from raw data via LLM' approach eliminates the need for any pre-existing knowledge graph or structured schemas, enabling the direct generation of meaningful data from the ground up, starting solely with company names. This positions our system as a vital tool for direct conversion of web-native, unstructured information into a usable structured format.

The demonstrated capability to identify and structure complex subsidiary relationships, especially for conglomerate companies (averaging 25.44 subsidiaries each), further validates the critical role of such automated systems. This level of detail in organizational complexity is challenging to ascertain manually and is crucial for comprehensive business intelligence, aligning with the complex relationship mapping endeavors highlighted in both EICopilot's graph-based exploration and the general need for semantic understanding in enterprise data analysis.

Compared to existing open-source frameworks, such as the LangChain AI *company-researcher*[12], which served as our foundational architecture, our system introduces substantial modifications. The most significant of these is the complete redesign of the information gathering mechanism. Rather than relying on prebuilt third-party solutions such as Tavily[13] or other search APIs, we implemented a fully custom web automation layer to extract domain-specific data with greater control and accuracy.

In addition to these infrastructural changes, we introduced targeted improvements in prompt engineering, a critical component in shaping the behavior and performance of language model-based agents (31; 32; 33; 34). While agent prompts were initially sourced from the LangChain repository, they were customized to align with our specific use case requirements. Additional enhancements include the integration of Pydantic-based schema validation, scalable bulk processing workflows, and a modular architecture designed for production-grade deployment. These improvements collectively enable our system to address the demands of enterprise-level applications, particularly within the Turkish-language market.

While our study demonstrates the technical feasibility of automated company data enrichment, it does not quantita-

---

[11]https://tinkerpop.apache.org/
[12]https://github.com/langchain-ai/company-researcher
[13]https://www.tavily.com

tively assess the downstream business impact of the enriched dataset on organizational decision-making processes. Future work should incorporate impact assessment methodologies to evaluate how enriched company profiles influence key business processes. Controlled studies comparing business outcomes using manually researched data versus automatically enriched profiles would provide valuable insights into the system's practical utility and identify areas where information quality most significantly influences business performance.

In addition, future research could explore integrating specialized extraction modules, similar to the annual report analysis method, to deepen the information extracted for specific document types (e.g., financial statements, legal filings) for companies where such structured reports are available. Additionally, enhancing the system's capabilities for automatic knowledge graph construction or populating existing graphs with the enriched data could further bridge the gap with knowledge graph-focused systems like EICopilot. Lastly, continued optimization of prompt engineering and LLM configurations, particularly for attributes with lower completion rates such as founder names, will be vital for maximizing the accuracy and completeness of the extracted business intelligence.

The system's performance characteristics would fundamentally change if *Gemini-2.0-flash*'s paid tier were utilized, primarily due to the elimination of rate limiting constraints. As demonstrated in our timing analysis in Table 2, our current implementation processes companies at an average of 169.78 seconds per company when constrained by free-tier limitations. In a paid-tier scenario, our architecture would enable parallel processing of companies, resulting in considerably shorter processing times. Given that our current implementation incorporates two 60-second sleep mechanisms to accommodate free-tier limitations, the actual computational overhead per company amounts to merely 49.78 seconds. This analysis suggests that with paid-tier access, our system could achieve substantially improved performance, processing entire datasets within a fixed time window rather than scaling linearly with dataset size.

## REFERENCES

[1] T. Grigalis and A. Cenys, "State-of-the-art web data extraction systems for online business intelligence," *Informacijos mokslai*, vol. 64, pp. 145–155, 01 2013.

[2] M. Arslan and C. Cruz, "Business-rag: Information extraction for business insights," in *Proceedings of the 21st International Conference on Smart Business Technologies - Volume 1: ICSBT*, INSTICC. SciTePress, 2024, pp. 88–94.

[3] S. Feuerriegel, J. Hartmann, C. Janiesch, and P. Zschech, "Generative ai," *Business Information Systems Engineering*, vol. 66, no. 1, pp. 111–126, feb 2024. [Online]. Available: https://doi.org/10.1007/s12599-023-00834-7

[4] J. S. Park, J. C. O'Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein, "Generative agents: Interactive simulacra of human behavior," 2023. [Online]. Available: https://arxiv.org/abs/2304.03442

[5] S. Hong, M. Zhuge, J. Chen, X. Zheng, Y. Cheng, C. Zhang, J. Wang, Z. Wang, S. K. S. Yau, Z. Lin, L. Zhou, C. Ran, L. Xiao, C. Wu, and J. Schmidhuber, "Metagpt: Meta programming for a multi-agent collaborative framework," 2024. [Online]. Available: https://arxiv.org/abs/2308.00352

[6] Q. Wu, G. Bansal, J. Zhang, Y. Wu, B. Li, E. Zhu, L. Jiang, X. Zhang, S. Zhang, J. Liu, A. H. Awadallah, R. W. White, D. Burger, and C. Wang, "Autogen: Enabling next-gen llm applications via multi-agent conversation," 2023. [Online]. Available: https://arxiv.org/abs/2308.08155

[7] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," 2023. [Online]. Available: https://arxiv.org/abs/1706.03762

[8] A. Radford and K. Narasimhan, "Improving language understanding by generative pre-training," 2018. [Online]. Available: https://api.semanticscholar.org/CorpusID:49313245

[9] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," 2020. [Online]. Available: https://arxiv.org/abs/2005.14165

[10] B. John, "Challenges and solutions in data integration for heterogeneous systems," 03 2025.

[11] H. Naveed, A. U. Khan, S. Qiu, M. Saqib, S. Anwar, M. Usman, N. Akhtar, N. Barnes, and A. Mian, "A comprehensive overview of large language models," 2024. [Online]. Available: https://arxiv.org/abs/2307.06435

[12] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J.-Y. Nie, and J.-R. Wen, "A survey of large language models," 2025. [Online]. Available: https://arxiv.org/abs/2303.18223

[13] C.-H. Jen, "Exploring construction of a company domain-specific knowledge graph from financial texts using hybrid information extraction," Dissertation, Stockholm, Sweden, 2021. [Online]. Available: https://urn.kb.se/resolve?urn=urn:nbn:se:kth:diva-291107

[14] K. Adnan and R. Akbar, "An analytical study of information extraction from unstructured and multidimensional big data," *Journal of Big Data*, vol. 6, no. 1, p. 91, 2019. [Online]. Available: https://doi.org/10.1186/s40537-019-0254-8

[15] J. Dagdelen, A. Dunn, S. Lee, N. Walker, A. S. Rosen, G. Ceder, K. A. Persson, and A. Jain, "Structured information extraction from scientific text with large language models," *Nature Communications*, vol. 15, no. 1, p. 1418, 2024. [Online]. Available: https://doi.org/10.1038/s41467-024-45563-x

[16] L. Ning, Z. Liang, Z. Jiang, H. Qu, Y. Ding, W. Fan, X. yong Wei, S. Lin, H. Liu, P. S. Yu, and Q. Li, "A survey of webagents: Towards next-generation ai agents for web automation with large foundation models," 2025. [Online]. Available: https://arxiv.org/abs/2503.23350

[17] Y. Snihur and J. Tarzijan, "Managing complexity in a multi-business-model organization," *Long Range Planning*, vol. 51, no. 1, pp. 50–63, 2018. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0024630117302893

[18] S. Yerva, Z. Miklos, and K. Aberer, "It was easy, when apples and blackberries were only fruits," vol. 1176, 10 2010.

[19] S. Datta, S. Bhattacharjee, and S. Das, "Clustering with missing features: a penalized dissimilarity measure based approach," *Machine Learning*, vol. 107, no. 12, pp. 1987–2025, Dec 2018. [Online]. Available: https://doi.org/10.1007/s10994-018-5722-4

[20] Lv, Zhepeng, "Business-oriented data management in conglomerates: insights from government data sharing," *E3S Web Conf.*, vol. 290, p. 02031, 2021. [Online]. Available: https://doi.org/10.1051/e3sconf/202129002031

[21] Y. Zhou, A. I. Muresanu, Z. Han, K. Paster, S. Pitis, H. Chan, and J. Ba, "Large language models are human-level prompt engineers," 2023. [Online]. Available: https://arxiv.org/abs/2211.01910

[22] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," 2023. [Online]. Available: https://arxiv.org/abs/2201.11903

[23] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, "Large language models are zero-shot reasoners," 2023. [Online]. Available: https://arxiv.org/abs/2205.11916

[24] D. Zhou, N. Schärli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, C. Cui, O. Bousquet, Q. Le, and E. Chi, "Least-to-most prompting enables complex reasoning in large language models," 2023. [Online]. Available: https://arxiv.org/abs/2205.10625

[25] B. Wang, S. Min, X. Deng, J. Shen, Y. Wu, L. Zettlemoyer, and H. Sun, "Towards understanding chain-of-thought prompting: An empirical study of what matters," 2023. [Online]. Available: https://arxiv.org/abs/2212.10001

[26] M. Kuutila, M. Mäntylä, and P. Raulamo-Jurvanen, "Benchmarking web-testing - selenium versus watir and the choice of programming language and browser," 2016. [Online]. Available: https://arxiv.org/abs/1611.00578

[27] A. Buscemi and D. Proverbio, "Chatgpt vs gemini vs llama on multilingual sentiment analysis," 2024. [Online]. Available: https://arxiv.org/abs/2402.01715

[28] H. Cai, X. Cai, J. Chang, S. Li, L. Yao, C. Wang, Z. Gao, H. Wang, Y. Li, M. Lin, S. Yang, J. Wang, M. Xu, J. Huang, X. Fang, J. Zhuang, Y. Yin, Y. Li, C. Chen, Z. Cheng, Z. Zhao, L. Zhang, and G. Ke, "Sciassess: Benchmarking llm proficiency in scientific literature analysis," 2024. [Online]. Available: https://arxiv.org/abs/2403.01976

[29] C. Viegas, R. Gheyi, and M. Ribeiro, "Assessing the capability of llms in solving poscomp questions," 2025. [Online]. Available: https://arxiv.org/abs/2505.20338

[30] Y. Yun, H. Ye, X. Li, R. Li, J. Deng, L. Li, and H. Xiong, "Eicopilot: Search and explore enterprise information over large-scale knowledge graphs with llm-driven agents," 2025. [Online]. Available: https://arxiv.org/abs/2501.13746

[31] P. Sahoo, A. K. Singh, S. Saha, V. Jain, S. Mondal, and A. Chadha, "A systematic survey of prompt engineering in large language models: Techniques and applications," 2025. [Online]. Available: https://arxiv.org/abs/2402.07927

[32] X. Amatriain, "Prompt design and engineering: Introduction and advanced methods," 2024. [Online]. Available: https://arxiv.org/abs/2401.14423

[33] S. Schulhoff, M. Ilie, N. Balepur, K. Kahadze, A. Liu, C. Si, Y. Li, A. Gupta, H. Han, S. Schulhoff, P. S. Dulepet, S. Vidyadhara, D. Ki, S. Agrawal, C. Pham, G. Kroiz, F. Li, H. Tao, A. Srivastava, H. D. Costa, S. Gupta, M. L. Rogers, I. Goncearenco, G. Sarli, I. Galynker, D. Peskoff, M. Carpuat, J. White, S. Anadkat, A. Hoyle, and P. Resnik, "The prompt report: A systematic survey of prompt engineering techniques," 2025. [Online]. Available: https://arxiv.org/abs/2406.06608

[34] J. White, Q. Fu, S. Hays, M. Sandborn, C. Olea, H. Gilbert, A. Elnashar, J. Spencer-Smith, and D. C.

Schmidt, "A prompt pattern catalog to enhance prompt engineering with chatgpt," 2023. [Online]. Available: https://arxiv.org/abs/2302.11382

• • •

**SAHIN UREGIL** is currently pursuing his B.Sc. degree in Computer Engineering from Galatasaray University, Istanbul, Turkey as of 2025. From July 2024 to June 2025, he worked as an R&D Intern at Kariyer.net, Istanbul, Turkey. He then continued as an AI Engineer Intern at TEB Arf from July 2025 to September 2025. His research interests include machine learning, deep learning, natural language processing, generative artificial intelligence, and applications of AI in finance.

**YAGIZ ERDEM** received his B.Sc. degree in Computer Engineering from Abdullah Gül University, Kayseri, Turkey, in 2025. He is currently pursuing his M.Sc. degree in Informatics at the Technical University of Munich, Germany. From July 2024 to September 2025, he worked as an R&D Intern at Kariyer.net, Istanbul, Turkey. His research interests include machine learning, deep learning, natural language processing, generative artificial intelligence, and computational drug discovery.

**SULTAN NEZIHE TURHAN** received her M.Sc. degree in Computational Sciences and Engineering from the Institute of Informatics at Istanbul Technical University in 2004, and her Ph.D. degree in Engineering Management from Marmara University in 2010. During her doctoral studies, she conducted research under the supervision of Prof. Hubert Kadima at École Internationale des Sciences du Traitement de l'Information (EISTI), focusing on healthcare supply chain process reengineering based on service-oriented architecture principles.

Between 1992 and 1998, she worked as a database administrator in various organizations. Since 1998, she has been serving as a lecturer at Galatasaray University, Faculty of Engineering and Technology, Department of Computer Engineering. Her research interests include Data Engineering, Database Management, Business Intelligence and Data Analytics, Educational Technologies, and Health Informatics.

In addition to teaching courses such as programming, database management systems, data analysis, business intelligence, and software engineering at different universities, she also provides programming education to individuals with special needs.

**GÜNCE KEZIBAN ORMAN** received her Ph.D. in Computer Engineering from INSA de Lyon. She is currently an Associate Professor in the Department of Computer Engineering at Galatasaray University, Istanbul, Turkey. Her research interests include graph neural networks, recommender systems, and complex network analysis, with a particular focus on accurate network modeling, community detection, and link prediction.

**DOGA YAGMUR YILMAZ** received her B.Sc. in Computer Engineering from Galatasaray University, Istanbul, Turkey, in 2023, and her B.Sc. in Industrial Engineering from Galatasaray University, Istanbul, Turkey, in 2024. She currently works as an R&D Engineer at Kariyer.net, Istanbul, Turkey. Her research focuses on developing AI-powered algorithms and digital recruitment solutions for candidate-employer matching and recommender systems, with a particular emphasis on suitability algorithms and data analysis.