



SCHOOL OF COMPUTATION,
INFORMATION AND TECHNOLOGY —
INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

**Reliable Multi-Agent Systems for
Automated Technical Data Acquisition and
Validation**

Yange Zheng





SCHOOL OF COMPUTATION,
INFORMATION AND TECHNOLOGY —
INFORMATICS

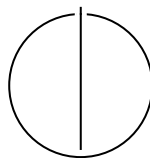
TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

**Reliable Multi-Agent Systems for
Automated Technical Data Acquisition and
Validation**

**Zuverlässige Multi-Agenten-Systeme für die
automatisierte technische Datenerfassung
und -validierung**

Author: Yange Zheng
Examiner: Prof. Dr. Viktor Leis
Advisor: Dr. Alexander Schiffmacher
Supervisor: Joe Yu
Submission Date: 01.07.2025



I confirm that this master's thesis is my own work and I have documented all sources and material used.

Munich, 01.07.2025

Yange Zheng

Acknowledgments

I owe special thanks to **Dr. Alexander Schiffmacher** (BMW AG), whose guides thesis from start to finish, with invaluable insights and support throughout the process.

I would also like to thank **Joe Yu** (BMW AG & TUM) for providing timely feedback and encouragement on key milestones.

Many thanks to **Stefan Merbele** (BMW AG) for his generous support with data access and infrastructure.

...

...

Finally, I am grateful to my family and friends for their unwavering support and encouragement during this journey.

Munich, July 2025

Yange Zheng

Abstract

In industrial procurement, engineers often work with **Quotation Analysis Forms (QAFs)** provided by suppliers—semi-structured documents containing part descriptions, electrical specifications, and inconsistent identifiers. While Manufacturer Part Numbers (MPNs) are typically present, they are often embedded in noisy text alongside other attributes, making reliable extraction and interpretation difficult. Automating the full pipeline—from QAF to datasheet to validated component specifications—requires coordinating diverse tools: language models, web search, document parsing, and schema validation. Rigid scripts and monolithic workflows tend to break in the face of data variation or missing fields. This thesis addresses that gap by designing a **modular multi-agent system** where each agent specializes in one task (e.g., MPN extraction, datasheet retrieval, VLM parsing), enabling flexible, fault-tolerant, and interpretable orchestration across real-world data acquisition scenarios.

Contents

Acknowledgments	iv
Abstract	v
1 Introduction	1
1.1 Background	1
1.2 Problem Statement	1
1.3 Research Objectives	2
1.4 Proposed Approach	2
1.5 Contribution	3
1.6 Thesis Structure	3
2 Related Work	4
2.1 Tool-Using Language-Model Agents	4
2.2 Vision-Language Models for Document Understanding	4
2.3 Multi-Agent Orchestration Frameworks	4
2.4 Document-AI in Procurement and Engineering	5
2.5 Research Gap	5
3 Problem Setting	6
3.1 Characteristics of Quotation Analysis Forms (QAFs)	6
3.2 Core Challenges	6
3.3 Formal Problem Definition	7
3.4 Automation Requirements	7
3.5 Mapping Challenges to Research Questions	8
4 System Architecture	9
4.1 High-Level Overview	9
4.2 Agent Responsibilities	9
4.3 Agent Coordination and Memory	10
4.4 Source Ranking Heuristic	11
4.5 Data Storage and Logging	11
4.6 Fault Tolerance	11

Contents

4.7 Security and Cost Considerations	11
4.8 Mapping to Requirements	11
5 Evaluation	13
6 Cost Model	14
7 Discussion	15
Abbreviations	16
List of Figures	17
List of Tables	18

1 Introduction

1.1 Background

The digitalisation of automotive supply chains has accelerated in recent years, yet many data-critical workflows still depend on semi-structured documents and manual interpretation. A prominent example at BMW is the *Quotation Analysis Form* (QAF), a file exchanged during supplier negotiations that lists part designations, preliminary prices, delivery conditions, and essential technical details. Although QAFs frequently include an identifier that resembles a Manufacturer Part Number (MPN), suppliers often embed the MPN inside longer descriptive strings or omit it altogether when quoting electrical characteristics. Consequently, engineers must open each QAF line, detect or correct the MPN, locate the corresponding datasheet, and then copy individual specifications—such as maximum voltage or operating temperature—into BMW’s internal cost and qualification systems. With several thousand parts assessed per programme phase, this *ad-hoc* process consumes hundreds of engineer-hours and remains prone to transcription errors.

1.2 Problem Statement

Preliminary attempts to automate QAF processing relied on regular expressions or supplier-specific parsing templates. Such approaches break when the supplier changes formatting, a datasheet link becomes unavailable, or unexpected language (e. g. German abbreviations) appears in the field. Modern language models (LLMs) and vision-language models (VLMs) promise the ability to reason over diverse text and visual layouts, yet their raw outputs can be inconsistent and lack provenance. A scalable solution must therefore *coordinate multiple specialised tools*, recover gracefully from partial failures, and guarantee that extracted numeric values are plausible and traceable to the source document.

1.3 Research Objectives

The central objective of this thesis is to design, implement, and evaluate a **reliable multi-agent system** that transforms noisy QAF inputs into validated, structured component specifications. To operationalise this goal, the work addresses three research questions:

- RQ1. How accurately can a language-model agent extract correct MPNs from semi-structured QAF strings?
- RQ2. Can cooperative retrieval agents—one using a distributor API, the other web search—achieve near-complete datasheet coverage while maintaining low latency?
- RQ3. What precision and recall can a vision-language parsing agent attain on key electrical specifications, and how do schema-based validation rules affect overall data quality?

1.4 Proposed Approach

To answer these questions, the thesis proposes a modular pipeline whose agents each address a single sub-task:

1. An **MPN Extraction Agent** uses GPT-4 with prompt engineering and regex post-filters to isolate standardised identifiers from QAF text.
2. Two **Retrieval Agents** query complementary sources: (i) the Octopart API for structured part metadata and (ii) a Google-backed web crawler prioritised by a cost-based heuristic that favours manufacturer domains.
3. A **Parsing Agent** leverages GPT-4 Vision to read PDF datasheets and return JSON fields—voltage, current, power rating, and temperature range—together with bounding-box provenance.
4. A **Validation Module** enforces JSON-Schema constraints and physical range rules and flags low-confidence extractions for manual review.
5. Finally, a **Cost-Prediction Model** (gradient-boosted regression) is trained on the extracted specifications, demonstrating the economic value of the pipeline’s output.

1.5 Contribution

The thesis offers four concrete contributions:

1. A fault-tolerant multi-agent framework that integrates LLMs, web search, distributor APIs, and VLM parsing into a single orchestrated workflow.
2. An implementation evaluated on **1 000** real BMW QAFs, including a manually annotated subset of 50–100 datasheets for ground-truth specification labels.
3. Empirical benchmarks reporting MPN extraction accuracy, datasheet retrieval success, field-level precision/recall, and end-to-end latency.
4. A downstream cost-prediction experiment showing that structured specs extracted by the system improve price-estimation accuracy, thus underlining direct business impact.

1.6 Thesis Structure

The remainder of this document is organised as follows: Section 2 surveys literature on language-model agents, vision-language models, and document AI in procurement. Section 3 formalises the challenges inherent to QAF processing. Section 4 details the proposed multi-agent architecture, while Section ?? describes implementation specifics, including prompts, API wrappers, and validation logic. Section 5 presents the experimental setup and quantitative results. Section 6 demonstrates the downstream cost-modelling application. Section 7 reflects on limitations and industrial implications, and Section ?? concludes the thesis with avenues for future work.

2 Related Work

This chapter surveys prior research relevant to the proposed multi-agent pipeline, grouped into four areas: (i) tool-using language-model agents, (ii) vision-language models for document understanding, (iii) multi-agent orchestration frameworks, and (iv) document-AI applications in procurement and engineering.

2.1 Tool-Using Language-Model Agents

Large language models have been extended with *tool use* capabilities to overcome context-length limits and provide up-to-date information. ReAct [yao2023react] interleaves chain-of-thought reasoning with executable actions, while the OpenAI function-calling interface popularised structured tool invocation in production systems. LangChain¹ offers an open-source abstraction layer for wrapping external tools (APIs, SQL, web search) behind LLM calls. AutoGPT and Voyager [wang2023voyager] push autonomy further but lack guarantees on reliability—highlighting the need for validation mechanisms such as those proposed in this thesis.

2.2 Vision-Language Models for Document Understanding

Document AI has progressed from OCR-centric pipelines to transformer-based models that combine layout and text. LayoutLMv3 [xuhuang2022layoutlmv3], DocPrompting [li2023docprompt], and StrucTexT [li2021structext] incorporate visual and spatial cues for table extraction. Recent work extends general VLMs (BLIP-2, GPT-4 Vision) to technical documents, but evaluation on complex engineering datasheets remains limited. Our parsing agent adopts GPT-4 Vision and focuses on reliability and schema-validated output rather than raw captioning.

2.3 Multi-Agent Orchestration Frameworks

Early agent systems in NLP used finite-state controllers; recent research revisits *LLM-driven* agents. MetaGPT [hong2023metagpt] assigns specialised roles (product manager,

¹<https://www.langchain.com>

engineer) to different LLM instances. CrewAI² and OpenAI’s “Assistants” API provide abstractions for role-based coordination. Most studies focus on qualitative demos; quantitative benchmarks on noisy industrial tasks, such as QAF processing, are scarce. Our work contributes empirical results on orchestration efficacy (success, latency) in a production-grade setting.

2.4 Document-AI in Procurement and Engineering

Prior industry case studies address invoice OCR, purchase-order matching, and part catalogue normalisation [chung2020invoiceai; koch2022bom]. Datasheet processing is typically rule-based or limited to text segments [rahman2021datasheet]. This thesis fills that gap and releases a benchmark on 1 000 QAFs to foster reproducible research.

2.5 Research Gap

- Existing LLM agents excel in controlled benchmarks but lack validation and provenance when applied to safety-critical engineering data.
- VLMs demonstrate impressive captioning accuracy yet are rarely evaluated on fine-grained numeric specs inside technical PDFs.
- No prior work jointly optimises extraction *and* downstream business value (e.g. cost prediction) within an end-to-end industrial pipeline.

These gaps motivate the multi-agent system and evaluation strategy proposed in the following chapters.

²<https://github.com/joaomdmoura/crewAI>

3 Problem Setting

3.1 Characteristics of Quotation Analysis Forms (QAFs)

A Quotation Analysis Form (QAF) is a semi-structured document exchanged during early supplier negotiations. At BMW, QAFs appear in multiple digital formats—CSV exports, Excel sheets, and supplier-specific PDF templates. Typical fields include:

- **Part Designation:** free-text description (e. g., “STM32F103C8 MCU 64 KB Flash TQFP48”),
- **Technical Specs:** one or more numeric attributes (voltage, current, temperature limits),
- **Price and MOQ:** preliminary unit price, minimum order quantity,
- **Manufacturer Part Number (MPN):** often concatenated with package or tolerance, occasionally missing or misspelled.

The absence of a universal schema means column order, header names, and punctuation vary widely across suppliers. Furthermore, multilingual notations (“Gehäuse”, “Temp. Bereich”) complicate rule-based parsing.

3.2 Core Challenges

C1: Inconsistent Identifier Extraction An MPN may be surrounded by package codes (“-TR”), voltage values, or marketing text, making pattern-matching unreliable. Example:

PI3USB9281ZLE_EX 20V Vbus Switch, TQFN-14, -40-85°C

C2: Heterogeneous Datasheet Sources Even after a correct MPN is found, the corresponding datasheet might be located on:

1. manufacturer domains (trusted),

2. distributor portals (moderately trusted),
3. third-party PDF mirrors (low trust, risk of 404 or captcha).

Selecting the “best” source requires ranking by reliability, download cost, and historical success.

C3: Complex Document Layouts Technical PDFs mix tables, figures, and footnotes. Numeric specs can appear in summary tables or deep in application sections. Standard OCR engines often split rows or mis-align columns, reducing extraction accuracy.

C4: Lack of Validation and Provenance Classic table-extraction pipelines output values without units or context. For safety-critical components, engineers must know *where* a value came from (page and bounding box) and whether it respects physical constraints (e. g., voltage > 0).

3.3 Formal Problem Definition

Let $q \in \mathcal{Q}$ denote a single QAF entry containing noisy text tokens. The task is to compute a validated specification record

$$\mathbf{s} = (\text{MPN}, \text{Voltage}_{\max}, \text{Current}_{\max}, \text{TempRange}, \dots)$$

such that:

1. **Correctness:** each numeric field equals the value printed in the authoritative datasheet.
2. **Completeness:** mandatory fields are present; optional fields are null if not in the document.
3. **Provenance:** every field carries evidence page, bbox.
4. **Latency:** average end-to-end processing time < 15 s per QAF.

3.4 Automation Requirements

Derived from the challenges and formal criteria, the target system must:

1. **Robust MPN Extraction** – tolerate noise, multi-language descriptors, and embedded package codes.

2. **Adaptive Retrieval** – combine distributor APIs and web search to maximise datasheet coverage while minimising failed downloads.
3. **Multi-Modal Parsing** – handle tables, prose, and regulatory footnotes using a vision-language model.
4. **Schema Validation** – enforce type, unit, and range checks; flag anomalies.
5. **Traceability** – store page coordinates and raw snippets for audit.
6. **Scalability** – process at least 3 000 QAF lines per day on commodity cloud hardware.

3.5 Mapping Challenges to Research Questions

Challenge	Requirement	Mapped RQ
C1	R1	RQ1
C2	R2	RQ2
C3	R3	RQ3
C4	R4, R5	RQ3 (validation aspect)
Performance	R6	All RQs (latency+throughput)

This mapping clarifies how each experimental evaluation in later chapters directly addresses an identified pain point in BMW’s procurement workflow.

4 System Architecture

This chapter describes the overall design of the proposed multi-agent pipeline, the responsibilities of each agent, their communication strategy, and the supporting infrastructure for storage, logging, and fault-recovery.

4.1 High-Level Overview

Figure 4.1 illustrates the data flow from a raw QAF line to a validated specification record and a downstream cost prediction. Each rectangular node represents an autonomous agent implemented as a function-calling wrapper around a large language model (LLM) or vision-language model (VLM). Solid arrows indicate the primary execution path, while dashed arrows represent feedback loops used for retries and validation.

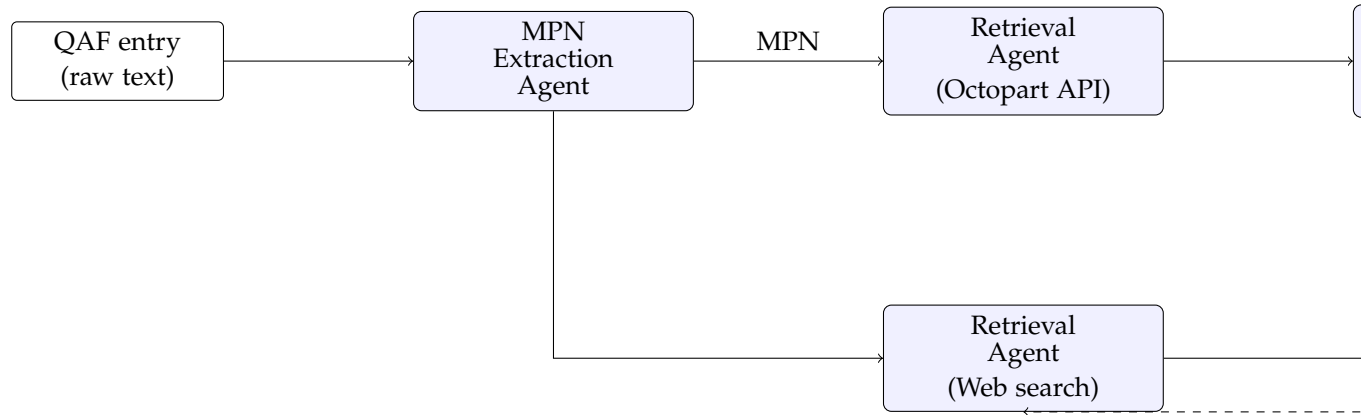


Figure 4.1: End-to-end architecture of the multi-agent pipeline.

4.2 Agent Responsibilities

MPN Extraction Agent Receives a raw QAF string and applies a few-shot GPT-4 prompt that asks the model to return a JSON object with the most likely manufacturer

part number. A lightweight post-processor removes package suffixes (e. g. “-TR”) and checks against a regex whitelist of vendor prefixes.

Retrieval Agent (API) Queries the Octopart REST API with the extracted MPN. If the API returns datasheets with a direct PDF link and the file size is within 50 kB–20 MB, the link is forwarded to the Parsing Agent.

Retrieval Agent (Web) Uses Playwright to perform a Google search of the form {MPN} datasheet filetype:pdf. The top-3 PDFs are downloaded for validation.

Parsing Agent Given a PDF path, calls GPT-4 Vision with a chain-of-thought template: first summarise the datasheet sections, then extract `max_voltage`, `max_current`, `power_rating`, `temp_range`, and return JSON plus page/bbox evidence.

Validation Module Executes three layers of checks:

1. JSON-Schema validation (types + units),
2. range rules (e. g. $0 < V_{max} < 1000V$),
3. duplicate consistency (table vs. paragraph vs. Octopart value).

Failures trigger a dashed feedback edge to the Web Retrieval Agent for retry.

Cost Model Trains an XGBoost regressor on validated specs + historical SAP prices to demonstrate the economic relevance of the structured output.

4.3 Agent Coordination and Memory

Agents communicate via the LangChain executor, which offers:

- **Function-calling interface:** each agent exposes a JSON schema; GPT-4 picks the correct function based on the task.
- **Shared Memory:** RedisJSON stores intermediate results (`mpn`, `pdf_path`, validation flags) accessible to all agents.
- **Retry Logic:** a decorator intercepts exceptions, waits exponentially (1 s, 2 s, 4 s), and reroutes control to the fallback agent if attempts > 3 .

4.4 Source Ranking Heuristic

For every candidate PDF URL u , the cost score is computed as

$$\text{cost}(u) = 0.2 \text{ domain}(u) + 0.4 \text{ pdfWeight}(u) + 0.4 \text{ failRate}(u), \quad (4.1)$$

where $\text{domain} = 0$ for manufacturer sites, 1 for trusted distributors, 2 otherwise; $\text{pdfWeight} = 0$ if the URL ends with `.pdf`, else 1; and failRate is an exponentially decayed average of previous download failures from the same host.

4.5 Data Storage and Logging

Validated JSON records are written to a PostgreSQL table with columns for `mpn`, `spec_json` (JSONB), `provenance` (JSONB), and an Azure Blob link to the PDF. Open-Telemetry traces (agent name, latency, token usage) are exported to Grafana for live monitoring.

4.6 Fault Tolerance

If both retrieval agents fail, the system stores a `status = "unresolved"` entry for manual triage. Parsing errors with low OCR confidence (< 0.7) raise a warning but still save partially validated fields, ensuring *at-least-once* data capture.

4.7 Security and Cost Considerations

All external calls pass through BMW's outbound proxy; API secrets are managed via Azure Key Vault. Estimated operating cost on an Azure D8s v5 VM plus one A10 GPU is ≈ 0.90 per 1 000 processed QAF lines, dominated by GPT-4 Vision tokens.

4.8 Mapping to Requirements

Table 4.1 summarises how architectural components satisfy the automation requirements defined in Chapter ??.

Together, these design choices create a fault-tolerant, interpretable, and extensible architecture capable of meeting BMW's data-quality and throughput requirements.

Table 4.1: Requirements coverage by architectural component.

Requirement	Fulfilled by
R1 Robust MPN extraction	MPN Extraction Agent (LLM + regex)
R2 Adaptive retrieval	API + Web Retrieval Agents, heuristic ranking
R3 Multi-modal parsing	VLM Parsing Agent (GPT-4 Vision)
R4 Schema validation	Validation Module (JSON Schema, range rules)
R5 Traceability	Provenance storage (page, bbox)
R6 Scalability	Redis memory, async Playwright, GPU offload

5 Evaluation

6 Cost Model

7 Discussion

Abbreviations

List of Figures

4.1	End-to-end architecture of the multi-agent pipeline.	9
-----	--	---

List of Tables

4.1	Requirements coverage by architectural component.	12
-----	---	----