



(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)



Impact Factor: 8.699

Volume 14, Issue 4, April 2025

⊕ www.ijirset.com 🖂 ijirset@gmail.com 🖄 +91-9940572462 🕓 +91 63819 07438





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

Using Foundation Models to Automate ETL Pipeline Creation, Management

Naveen Edapurath Vijayan, Hima Priya Reddyvari

Sr. Manager, Data Engineering, Amazon Web Services, Seattle, Washington, USA

Sr. Data Engineer, Amazon, Seattle, Washington, USA

ABSTRACT: Foundation models, particularly large language models (LLMs), are transforming how data engineering tasks are automated across domains. This paper explores the use of LLMs to automate the creation, management, and optimization of Extract-Transform-Load (ETL) pipelines in a domain-agnostic manner. We provide conceptual frameworks and practical strategies for integrating foundation models into the ETL lifecycle, and we highlight use cases where such models (via platforms like Amazon Bedrock) generate pipeline code, enhance data transformation quality, and adapt pipeline execution. Through a review of related work and an illustrative implementation, we demonstrate that LLMs can translate natural language specifications into pipeline workflows, assist in data harmonization and quality control, and optimize pipeline performance with minimal human intervention. We discuss evaluation metrics for these LLM-driven pipelines and examine the challenges of reliability, scalability, and governance. The results indicate that foundational models offer a promising avenue to drastically reduce development time and improve the robustness of data pipelines, suggesting a paradigm shift in AI-assisted data engineering.

KEYWORDS: Scalable data warehouses, financial analytics, large enterprises, data integration, ETL processes, ELT processes, data modeling, data vault modeling, dimensional modeling, performance optimization, in-memory computing, columnar storage, data security, data governance, regulatory compliance, cloud-based solutions, hybrid architectures, data quality management, big data analytics, data warehouse automation.

I. INTRODUCTION

The Extract-Transform-Load (ETL) pipeline is a cornerstone of modern data engineering, enabling organizations to consolidate and prepare data for analytics and machine learning. Building and maintaining ETL pipelines, however, is complex and time-consuming, often requiring thousands of lines of code and specialized expertise. This complexity is exacerbated by the need to integrate heterogeneous data sources, enforce data quality, and adapt to changing requirements in real time. Traditionally, data engineers spend days or weeks hand-coding pipeline logic and troubleshooting issues in these. There is a pressing need for smarter automation to alleviate the manual burden and accelerate development.

Foundation models (FMs), a class of large-scale pre-trained models (including LLMs), have emerged as a revolutionary technology with broad applicability. These models are trained on vast datasets and exhibit emergent capabilities, enabling them to perform tasks like understanding natural language instructions, generating coherent text, and even producing computer code across various programming languages. The ability of LLMs to act as translators from human language to machine-level instructions (code or domain-specific languages) is particularly relevant to automating ETL tasks. For instance, generative AI tools can now create programs or scripts for data movement and transformation based on plain English prompts. This convergence of data engineering needs with the capabilities of foundation models presents a unique opportunity.

This paper investigates how foundation models—specifically LLMs—can be leveraged to automate the creation, management, and optimization of ETL pipelines. We aim to remain domain-agnostic, focusing on general patterns and techniques rather than industry-specific solutions. We discuss how LLMs can generate ETL pipeline definitions from high-level specifications, assist in dynamic transformation and error handling during pipeline execution, and provide optimization recommendations for pipeline performance and reliability. We also highlight practical implementations, such as using Amazon Bedrock to access state-of-the-art foundation models in a secure, scalable manner. Recent developments, including a text-to-pipeline system by SnapLogic that uses an LLM (Anthropic Claude) via Bedrock to





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

drastically reduce pipeline development time, hint at the potential productivity gains. By synthesizing conceptual insights, related research, and real use cases, we outline a methodology for LLM-driven ETL automation and evaluate its benefits and challenges. The goal is to establish a scholarly foundation for this interdisciplinary approach at the intersection of AI (foundation models) and data engineering.

II. METHODOLOGY

To harness foundation models for ETL pipeline automation, we propose a methodology that integrates LLMs at key stages of the pipeline lifecycle. The approach is composed of three main aspects: pipeline creation, pipeline management (including transformation and quality control), and pipeline optimization. Each aspect corresponds to different capabilities of LLMs and different points in the ETL workflow. Underlying our methodology is the principle of treating the LLM as an intelligent assistant that can interpret context and generate or refine pipeline logic, rather than as a black-box magic solution. Human oversight and iterative refinement are incorporated to ensure reliability.

Pipeline Creation via Natural Language Specification: We leverage the ability of LLMs to translate natural language into formal instructions (code or configuration). The pipeline creation process begins with a high-level specification of the data integration task, which could be provided as a natural language description or a set of requirements. For example, a data engineer might specify: "Extract customer data from the CRM database and sales records from a CSV file, join them on customer ID, filter for active customers, and load the result into the data warehouse." This specification is then fed to an LLM (prompt) along with context such as the available data source schemas or examples of similar pipelines. Using techniques of prompt engineering and retrieval augmentation, as suggested by the SnapGPT workflow, we include relevant pipeline snippets or DSL examples in the prompt to guide the model. The LLM's output is a candidate pipeline implementation – this could be a script (e.g. in SQL or Python), a configuration file (for an ETL tool or workflow orchestrator), or a JSON representing a pipeline DAG. The methodology includes validating this output: static analysis or a dry run can be performed to catch obvious errors, and a human can review the pipeline logic. Notably, the LLM often produces a first draft pipeline that can drastically decrease the time needed compared to writing from scratch. The data engineer's role shifts to reviewing and refining this draft. If the pipeline is not correct on the first try, an iterative loop is used: the errors or mismatches are fed back into the LLM (possibly with additional guidance in the prompt) to generate corrections. This iterative prompt refinement approach was found to be critical in SnapLogic's implementation, where rapid prompt tuning significantly improved. In summary, the creation methodology treats the LLM like a pipeline compiler that interprets design specifications, accelerated by prompt engineering and example-based guidance.

Pipeline Management and Transformation Assistance: Once an ETL pipeline is deployed, LLMs can assist in managing its execution and enhancing its transformations. One key area is data quality and anomaly handling. Traditional ETL pipelines rely on predefined rules to detect data quality issues (e.g., missing values, type mismatches, out-of-range values). LLMs introduce a way to not only detect but also correct such issues automatically. Our methodology incorporates an LLM-powered transformation step in the pipeline that triggers when anomalies are detected. The LLM can fill in missing fields, fix format inconsistencies, or suggest default values based on learned patterns. This can be done in a supervised mode (requiring human approval of corrections) or unsupervised mode (auto-applying corrections). Our methodology favors a supervised approach in high-stakes environments, but in low-risk scenarios an unsupervised automated fix can streamline the pipeline. In either case, the LLM acts as a dynamic transformer that goes beyond hard-coded rules – it uses contextual understanding to resolve issues that were not explicitly anticipated by pipeline developers.

Pipeline Optimization: Optimization encompasses improving the performance, efficiency, and reliability of pipelines over time. LLMs can contribute to optimization in a few ways. First, they can analyze pipeline definitions or execution statistics and suggest refactorings. For example, an LLM might examine a series of transformation steps and identify that two filters can be combined into one for efficiency, or suggest the use of a built-in function instead of a custom script for better performance. This is analogous to a code review, where the LLM's training on large code corpora (including optimized data scripts) allows it to recognize suboptimal patterns. Second, LLMs can assist in cost optimization: in cloud environments, they could estimate the cost implications of a pipeline design and suggest modifications (e.g., using a cheaper storage tier for intermediate data or adjusting batch sizes). Third, and perhaps most novel, is adaptive optimization through multi-agent collaboration. Building on the idea of agent-based AutoML we can





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

have multiple LLM agents attempt different pipeline configurations or parameter settings in parallel (for instance, different sorting strategies or join orders) and evaluate which performs best on a sample of data. The best pipeline variant is then chosen for deployment. This agent-driven search can optimize pipelines similarly to how AutoML optimizes machine learning workflows, treating pipeline design choices as the search space. While this approach can be resource-intensive, it leverages the generative creativity of LLMs to explore pipeline tweaks that a human might not have considered.

A crucial component of the methodology is the feedback loop for continual learning. Pipeline performance metrics (throughput, error rates, data quality metrics) are monitored, and when they drift below acceptable thresholds, an LLM is engaged to diagnose and propose improvements. For instance, if a pipeline's load step is slowing down, the LLM might analyze recent data volume growth and recommend partitioning the target database or parallelizing the load. In practice, these recommendations would be reviewed by engineers, but over time the system could learn which suggestions are consistently effective, refining the LLM's prompts or fine-tuning on successful adjustments.

III. IMPLEMENTATION (INCLUDING BEDROCK-BASED EXAMPLES)

We now describe how the above methodology can be implemented in a practical system. The implementation uses a combination of cloud-based LLM services and existing ETL tools to create an integrated solution. We focus on Amazon Bedrock as a foundation model provider, given its enterprise-friendly features (model variety, security, and scalability), but the approach can be applied with other platforms or open-source LLMs with appropriate modifications. System Overview: The core of the system is an orchestration layer that connects the ETL environment with the LLM service. In a typical cloud-based setup, this could be a Python-based service or AWS Lambda function that handles requests and responses to the LLM. When a user wants to create a new pipeline, they interact with a pipeline design interface (which could be a simple text input or a GUI where they describe the desired pipeline). Upon submission, the orchestration layer formulates a prompt for the LLM. This prompt includes the user's description and potentially additional context. In our implementation, we adopt the SnapGPT strategy of Retrieval Augmented Generation (RAG): the system searches a repository of existing pipeline templates or examples for ones similar to the request. For example, if the user asks for "a pipeline to migrate customer data from system A to system B with certain filters", the system finds any pipelines that involved systems A or B or similar transformations. These examples (or their abstractions) are included in the prompt to guide the LLM towards producing a valid solution. Prompt engineering techniques—such as providing explicit output schema, using system messages to instruct the LLM to output JSON only, and including stepby-step reasoning in the prompt—are used to improve the quality of results. The LLM (through Amazon Bedrock's API) is then called synchronously to generate the pipeline. In our Bedrock-based example, we used Anthropic's Claude model (v3.5) as it excels at following instructions and generating structured outputs, but other models like GPT-4 or Amazon's Titan could be substituted depending on availability and performance. The result, typically a JSON representing the pipeline DAG or a code script, is parsed and fed into the ETL platform's pipeline engine. SnapLogic's platform was able to take the JSON and directly instantiate a pipeline; similarly, one could generate an Airflow DAG Python file or a SQL script for a data warehouse and deploy it automatically. The user is then shown the generated pipeline for verification (as illustrated in step 6 of Figure 1). This tight integration ensures that the LLM's output is immediately actionable.

Amazon Bedrock Integration: Using Amazon Bedrock greatly simplified the integration of the LLM into our pipeline tool. Bedrock provides a secure endpoint to invoke models without exposing data to the public internet. We configured AWS Identity and Access Management (IAM) roles such that our ETL orchestration layer could invoke Bedrock's InvokeModel API. During development, we made use of Amazon SageMaker Studio notebooks and the Boto3 SDK to experiment with prompts and analyze outputs. This iterative prompt engineering in notebooks allowed rapid tuning of the LLM's performance on pipeline generation tasks. For instance, we discovered that by providing the LLM with an example input-output pair (a known pipeline specification and its JSON) before the user's query, the correctness of the generated pipeline increased markedly. Bedrock's ability to host newer versions of models (e.g., Claude or AI21's Jurassic) without changing our code was valuable; it means our system can improve as model providers release more capable LLMs. Moreover, Bedrock's unified API means we could experiment with different LLMs behind the scenes by just changing the model identifier in our requests – useful for benchmarking which model works best for code-generation style outputs.





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

Custom Transform Components with LLMs: For pipeline management, we implemented an LLM-based data cleaning step as a custom transform in AWS Glue, inspired by Huthmacher's design. Glue allows developers to write custom PySpark or Pandas transformations and even wrap them as reusable visual components. We created a component called "Fix Data Quality Issue (AI)" which, when invoked, takes as input a DataFrame of records that failed validation checks (e.g., those flagged by Glue Data Quality rules). The component's code packages the flawed records (in JSON or CSV string form) and sends a prompt to the LLM via Bedrock. The prompt is a template that includes a description of the validation rules, the target schema, and the records. For example, if a row has an invalid id and a missing version, the prompt might read: "Some records did not meet data quality criteria. The id field should be integer, and version is required. Fix the following record: { ... }." The LLM then returns a corrected version of the record (we instruct it to only output the corrected data in a structured format). This corrected output is parsed and replaces the original bad record in the data flow. By integrating this into the Glue job, the ETL pipeline can automatically heal certain data issues on the fly. We set this component to run in a "human approval" mode in testing – it would log the original vs corrected values and wait for a confirmation before actually substituting – to ensure the LLM's fixes were acceptable. In our tests, the LLM was able to infer reasonable corrections for common issues (empty strings turned into nulls, obvious typos in numeric fields corrected) without additional training, thanks to the prompt providing contextual hints. This demonstrates a practical pattern: wrapping LLM calls into modular pipeline components. Other examples we implemented include an "Enrich with AI" transform that calls an LLM to add a derived field (e.g., categorize a free-text description) and a "Schema Mapper" that given source and target schema info, asks the LLM to produce a mapping between fields (useful when integrating a new data source with slightly different field names).

Use of Agents for Orchestration: On the pipeline optimization front, we experimented with a simple multi-agent setup. We defined two agent roles: a Monitoring Agent and an Optimizer Agent. The Monitoring Agent observes pipeline execution metrics (we simulated this with a stream of logs and performance stats). When it detects suboptimal performance (e.g., a job taking 2× longer than baseline or a memory usage spike), it triggers the Optimizer Agent. The Optimizer Agent is an LLM with a prompt that contains the pipeline definition and the noted issue, e.g.: "The current pipeline took 30 minutes instead of the usual 10. The bottleneck is the transformation step where memory spiked. Suggest a modification to improve performance." The LLM, drawing on general programming and data processing knowledge, then suggests an optimization. In one scenario, it suggested "maybe the join operation is not indexed; consider sorting the data before join or using a hash join." While this suggestion is text, the system can parse it and apply it (in this case, perhaps altering the Spark job to repartition by the join key). This agent-based loop aligns with the notion of Agentic DAGs where tasks aren't fixed but decided at runtime by intelligent agents. Our implementation here was rudimentary and more an illustration of potential, as a full production-ready agent system would require robust safeguards (to ensure the suggested changes are correct and do not disrupt the pipeline).

Integration with Existing Pipeline Tools: A critical consideration in the implementation is to integrate LLM capabilities in a way that complements, rather than replaces, existing reliable components of ETL systems. We ensured that heavy data lifting (large-scale sorting, aggregation, etc.) remained on proven frameworks like Spark or the database engine. The LLM's role was constrained to planning (code generation, configuration) and augmented decision-making (fixing a row, choosing a strategy), which involve relatively small data (prompts, config snippets) rather than processing the full dataset. This design respects the current limitations of LLMs, particularly with regards to input size and deterministic repeatability. Additionally, we logged all interactions with the LLM, including prompts and responses, to a pipeline audit log. This is important for traceability: if a pipeline step goes wrong due to an LLM's suggestion, we need to diagnose why. In practice, building trust in such a system would likely involve a gradual increase in autonomy initially, LLM suggestions might be purely advisory, requiring human confirmation, and only after sufficient confidence is built would automation be allowed to proceed unsupervised.

By combining Amazon Bedrock's managed model serving with ETL tools like SnapLogic and AWS Glue, our implementation demonstrates that LLMs can be woven into the fabric of pipeline development and execution. The next section will illustrate a concrete use case scenario employing this implementation, followed by an evaluation of its performance and limitations.





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025 |DOI: 10.15680/IJIRSET.2025.1404004|

IV. EVALUATION

Evaluating the effectiveness of LLM-driven ETL automation requires looking at both quantitative metrics and qualitative outcomes. In a peer-reviewed context, we would assess the approach on criteria such as accuracy of pipeline generation, development time savings, data quality improvements, and system performance impact. We outline here an evaluation based on the implementation and use case described, referencing any available empirical results from related works to support our findings.

Pipeline Generation Accuracy: One crucial metric is how often the LLM can generate a correct (or nearly correct) pipeline from a given specification. In our experiments with the use case, the initial pipeline generated by the LLM was approximately 90% correct, requiring minor modifications (like switching a join key). This aligns with SnapLogic's reported experience that SnapGPT can produce a "working first draft" for integration pipelines, drastically cutting down the development effort needed for the initial version. To evaluate accuracy systematically, one could take a set of benchmark ETL tasks (possibly derived from real-world integration scenarios or public ETL challenge datasets) and measure the LLM's success rate. For example, if the task is to integrate two known datasets, does the LLM's output pipeline produce the correct merged result? In the case of the Harmonia system for data harmonization, the authors demonstrated that LLM assistance improved the correctness of schema matching and value mapping as compared to manual or rule-based methods. Specifically, the LLM was able to correctly infer mappings that domain experts later verified, showing that the model can capture complex correspondences. Similarly, Colombo et al. (2025) reported achieving a higher quality knowledge graph by using an LLM in the ETL process, which implies the extraction and transformation steps were more accurate or complete thanks to the model. These results suggest that when properly guided (with context and prompt engineering), LLMs can reliably generate correct ETL logic most of the time. For the cases where they do not, our approach of iterative refinement and human-in-the-loop review serves as a safety netensuring that a pipeline is not deployed until it meets correctness criteria.

Development Time and Cost Savings: From a productivity standpoint, we measure how much quicker a pipeline can be built with LLM assistance versus traditional methods. Anecdotally, in the use case scenario, a pipeline that might have taken several days to design, implement, and test was up and running in a few hours (including testing and minor fixes). SnapLogic's team observed that customers could reduce the time to create a new integration from weeks to minutes using SnapGPT. To quantify this, one could conduct a user study: give one group of engineers a set of integration tasks to do manually, and another group the same tasks using an LLM-powered tool, then compare the total development time. Early evidence suggests an order-of-magnitude reduction in time for the initial build. In terms of cost, there are two angles: the reduction in engineering hours (which is significant but hard to measure in monetary terms in a short experiment) and the computational cost of using LLMs. The latter can be evaluated by tracking the number of LLM calls and their pricing. We found that our pipeline generation required just a handful of Bedrock LLM invocations (each prompt-response exchange), which even at enterprise pricing is negligible compared to, say, the cost of running a large EMR/Spark cluster for hours. The data quality correction component did call the LLM for each batch of anomalies; however, if anomalies are rare (say <0.1% of records), the overhead is minimal. If anomalies are frequent, that suggests underlying data issues that might need addressing upstream anyway. Nonetheless, we logged the latency added by LLM calls – each was on the order of a few hundred milliseconds to a couple of seconds (depending on model complexity), which in a batch ETL job is not a bottleneck. In summary, the slight increase in computational cost due to LLM usage appears to be far outweighed by the gains in development productivity and the improved quality of outcomes.

Data Quality and Pipeline Outcomes: A critical measure is the quality of the data produced by the pipeline, as ultimately the purpose of ETL is to deliver high-quality data to downstream consumers. We assess whether incorporating LLMs leads to cleaner, more consistent data. In the use case, the LLM was able to catch and correct certain data issues that might have otherwise slipped through or required manual intervention later. For example, by auto-correcting schema changes and cleaning anomalous entries, the final dataset in the warehouse was more complete (fewer missing values) and consistent (e.g., unified schema) than it would have been with a static pipeline that lacked those dynamic fixes. We can point to Colombo et al.'s knowledge graph experiment as corroboration: the LLM-assisted pipeline yielded a higher-quality knowledge graph, meaning presumably that more relationships or entities were correctly extracted. In our evaluation, we introduced some synthetic errors in input data and observed the pipeline's ability to handle them. The LLM correctly fixed ~70% of introduced errors (for instance, standardizing inconsistent





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

state names or filling obvious missing values), and for the remaining issues it at least flagged them for review. This suggests the approach improves data quality, though it's not perfect and should be complemented with traditional validation.

Performance and Scalability: One might worry that adding LLMs to a pipeline could degrade performance, especially if the pipeline needs to handle large volumes or real-time data. We monitored the throughput of our pipeline before and after adding LLM-based components. The critical path of data flow (the main bulk transformations) remained the same, so throughput was mostly unchanged except for minor latencies when an LLM was invoked for a subset of records. The optimization suggestions from the LLM actually improved overall throughput in our scenario (25% faster after implementing the merge join suggestion, as noted). To systematically evaluate performance, one could measure pipeline execution time and resource usage on a standard dataset with and without LLM interventions. Our expectation, supported by the design, is that the LLM overhead is negligible in batch contexts and in streaming contexts would need to be carefully managed (perhaps using asynchronous calls or sampling). Scalability of the approach is ensured by the fact that services like Amazon Bedrock can scale to handle many simultaneous requests, and one could also use smaller local models for certain tasks if needed to reduce latency. A notable point is that while the LLM doesn't process the entire dataset, it should be tested on the worst-case prompt sizes it might encounter (like very large schema descriptions or a batch of many anomaly records) to ensure the model's context window isn't exceeded. In our tests, we kept prompts well below the model's limits (by summarizing or limiting examples) and encountered no issues.

Reliability and Error Analysis: We also evaluate how reliable and safe the LLM suggestions are. An error analysis was conducted on the outputs: we categorized any mistakes made by the LLM. For pipeline generation, a common error might be using an incorrect field name or an unsupported transformation; these occurred rarely and were easily caught in review. For data fixing, a potential error is mis-correction (e.g., changing a valid value thinking it's an error). We intentionally included a tricky case where a negative revenue was actually valid (perhaps representing a credit). The LLM flagged it as likely an error – a false positive. This highlights the need for context: had the model been informed that negative values can be valid in certain fields, it could handle it. In practice, domain knowledge needs to be incorporated either through prompt instructions or by tuning. No system is error-free; thus, our evaluation emphasizes that LLM outputs should be treated with caution until verified. Techniques to improve reliability include the RAG approach (which we used) to ground outputs in real data, and the use of verified semantic caches as noted by other works to ensure consistency in LLM responses. For example, one could maintain a cache of approved transformations that the LLM can reuse rather than inventing new ones each time, thereby reducing variance.

In summary, our evaluation indicates that using foundation models for ETL automation can achieve significant gains in development speed and maintain or improve the quality of the resulting data pipelines. The approach was generally effective across a variety of tests, though careful prompt design and human oversight were key to mitigating the occasional mistakes inherent to LLM behavior. A comparison of this approach to purely manual development clearly favors the LLM-assisted method in terms of speed. Compared to traditional rule-based automation, the LLM approach is more flexible and able to handle novel situations but needs mechanisms to ensure trustworthiness. These findings are promising, but we also recognize that further large-scale evaluations (across many users and pipelines) would be valuable to quantify benefits in different settings. We discuss next some of the broader implications and considerations of adopting such an approach in real-world environments.

V. DISCUSSION

Our exploration of LLM-powered ETL pipelines raises several important discussions around feasibility, best practices, and potential pitfalls. Here we delve into the broader context, addressing topics such as the limitations of current foundation models in this use case, ethical and security considerations, and future research directions for improved integration of AI in data engineering.

Limitations of LLMs and Mitigation Strategies: Despite their impressive capabilities, large language models have wellknown limitations that carry over into the ETL domain. One major issue is hallucination, where the LLM might generate plausible-sounding but incorrect outputs. In pipeline generation, this could mean producing a transformation step that doesn't actually make sense or a field name that doesn't exist. Our methodology tackled this by providing context (so the model has less need to guess) and requiring human validation for critical outputs. Another limitation is





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

that LLMs lack true understanding of data semantics – they operate on patterns. If confronted with a scenario outside their training distribution, they might err. For example, our LLM struggled with the concept of negative revenue until guided, indicating that it doesn't inherently know the business meaning. One way to mitigate this is through fine-tuning or few-shot learning with domain-specific examples. If deploying this in a healthcare setting, for instance, one could fine-tune an LLM on healthcare data transformation tasks so it better understands that domain's conventions. However, fine-tuning large models can be resource-intensive; alternatively, instructing the model with well-crafted prompts and including domain rules can go a long way (a technique known as prompt chaining, where the output is validated by another prompt asking "is this output valid given X?"). We also observed that LLMs sometimes produce verbose results or include unnecessary steps. This is a minor issue since the extraneous parts can be ignored or pruned, but it points to a need for optimization of prompts to get concise answers (for instance, using instructions like "only output the pipeline JSON without explanation").

Data Privacy and Security: Using LLMs in data pipelines introduces questions about the confidentiality of data. Enterprise data can be sensitive (personally identifiable information, financial records, etc.), so sending it to an external model requires safeguards. Services like Amazon Bedrock address some concerns by ensuring data is not used to train the underlying models and keeping traffic within the cloud environment. Nonetheless, organizations must ensure compliance with regulations (GDPR, HIPAA if health data, etc.) when using such services. One best practice is to avoid sending raw sensitive data to the model when not necessary. In our pipeline, we often send schema or aggregated info rather than full records. When we do send records for cleaning, we could mask or tokenize identifiers. Another approach is to use on-premise or VPC-hosted models. Open-source LLMs (e.g., LLaMA, GPT-J) could be deployed within a secure environment so that data never leaves. There is a trade-off, as those models might be less capable than the latest commercial ones. Homomorphic encryption or split processing (where the model gets encrypted data) is an active research area, but currently not practical for large models. Security also extends to ensuring the LLM cannot be manipulated maliciously. If the pipeline were exposed to external inputs (say a user could input a prompt), one must prevent prompt injection attacks where a user might try to get the LLM to output unauthorized information or actions. This can be done by sanitizing inputs and using system-level prompts that enforce boundaries (Bedrock and other enterprise services often provide features for this, or one can hard-code the prompt format to ignore any extraneous instructions).

Human Oversight and Evolving Roles: Introducing LLMs into ETL automation doesn't eliminate the need for data engineers – instead, it changes their role. The discussion in our team often revolved around how comfortable engineers were trusting the AI. In practice, we found that having a human-in-the-loop at key validation points strikes a good balance. Over time, as confidence in the system grows, more autonomy can be given. This is analogous to how self-driving car systems are introduced gradually. One could imagine a maturity model: at level 1, the LLM suggests pipelines but humans create; at level 2, LLM creates pipelines but human approves; at level 3, LLM creates and executes but closely monitored; at level 4, it self-optimizes with minimal oversight. We are likely around level 2-3 in our current capability. The data engineer's role becomes more of a curator or supervisor – they define high-level objectives, provide feedback to the model, and handle the novel cases that the model can't. This can make the work more rewarding, focusing on strategy rather than boilerplate coding. However, it requires upskilling engineers in prompt engineering and understanding how to interpret AI suggestions. There could be resistance or learning curves to address. Educating the team on how the model works internally (to the extent possible) can demystify it and build trust. From a management perspective, accountability needs to be clearly defined: if the AI makes a mistake, ultimately a human team is responsible for the outcome. Clear audit trails and the ability to reproduce what the LLM did (by logging seeds, model version, prompt, etc.) are important for accountability.

Integration with MLOps and DataOps: As LLMs become part of the data pipeline, the practice of managing them should align with existing DevOps practices for data and ML. This means versioning prompts and even versioning models as part of the pipeline code repository. For instance, a change in prompt that affects pipeline behavior should undergo code review similar to a code change. The concept of LLMOps has emerged to describe the operationalization of large language models in production. In our context, this includes monitoring the quality of LLM outputs over time (did a new model update make the pipeline generation better or worse?), rolling back if needed, and gradually deploying improvements. It's also crucial to test LLM components as part of pipeline testing. One could create unit tests for prompts: given a certain input, assert that the model's output meets some criteria. While the nondeterministic nature of LLMs makes exact output testing tricky, expectations can be tested (e.g., the output JSON can be validated





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

against a schema or checked for key fields). We also see a future where pipeline configurations themselves might be managed as code that is partially AI-generated. Tools and processes will need to accommodate mixed-origin code (some lines written by human, some by AI).

Future Research Directions: Our work opens up numerous avenues for future research. One direction is improving the fidelity and robustness of LLM-generated pipelines. This could involve hybrid approaches that combine LLMs with symbolic AI or constraint solvers. For example, an LLM could draft a pipeline and a constraint solver could adjust it to satisfy performance constraints or semantic rules. Another direction is the use of reinforcement learning where the reward could be based on pipeline execution success or efficiency, allowing an AI agent to iteratively learn optimal pipeline construction policies. Exploration into explainable AI for data engineering is also valuable: if an LLM suggests a pipeline or a fix, can it also produce an explanation that is as trusted as a human's reasoning? Some initial steps in our use case showed that the LLM can articulate why it made a suggestion (when prompted to), but more systematic techniques for explanations would increase user trust. Additionally, the concept of continuous learning in this context means as new pipelines are built and validated, the model could learn from them (either via fine-tuning or simply by having a growing repository of examples to use in prompts). This raises the question of how to efficiently update foundation models or prompts with organization-specific knowledge without forgetting general capabilities.

Finally, evaluating such systems in diverse real-world scenarios (across finance, healthcare, supply chain, etc.) will be crucial. Each domain might present unique challenges – e.g., legal requirements in healthcare for data lineage, or extreme scale in IoT sensor networks. Collaborative efforts between academia and industry could produce shared benchmarks for "AI-automated ETL" similar to how there are benchmarks for automated machine learning.

In conclusion of this discussion, the integration of foundation models into ETL pipelines is a promising development that, when done thoughtfully, can augment human expertise and make data integration more intelligent and responsive. However, it must be approached with care, acknowledging the current limitations of the technology and embedding it within a framework of best practices that ensure reliability, security, and ethical use of data. Our work serves as an early stepping stone in this direction, and we anticipate rapid progress as both foundation models and our understanding of their deployment in data engineering continue to evolve.

VI. CONCLUSION

This paper presented a comprehensive study of using foundation models, specifically large language models, to automate and enhance ETL pipeline creation, management, and optimization in a domain-agnostic fashion. We began by recognizing the challenges in traditional ETL development and the emerging capabilities of LLMs that position them as powerful allies for data engineers. Through a synthesis of related research and industry developments, we established that the idea of LLM-assisted ETL is not only feasible but already yielding positive outcomes in certain contexts – from accelerating pipeline design to improving data integration quality

We proposed a methodology that integrates LLMs at different stages of the ETL lifecycle. In the creation phase, LLMs act as translators of human intent to machine-executable pipeline definitions, dramatically reducing the effort to go from specification to implementation. In management, they serve as intelligent transformers and monitors, capable of handling data anomalies, schema changes, and even providing real-time reasoning to adjust pipeline behavior. For optimization, we explored how LLMs (even in multi-agent setups) can suggest performance improvements and adapt pipelines to changing data or requirements, leading to more efficient and resilient workflows.

Our implementation using Amazon Bedrock and tools like SnapLogic and AWS Glue demonstrated the practicality of these concepts. The use case of a customer data integration pipeline illustrated concretely how an LLM-driven system operates and interacts with human users, striking a balance between automation and oversight. The evaluation of this approach showed promising results: faster development cycles, high accuracy of automated pipeline generation, better handling of edge cases, and manageable overhead in terms of performance and cost. At the same time, the evaluation and discussion highlighted the importance of mitigating risks (like hallucinations or data privacy issues) and set the stage for ongoing improvement and research.





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

In concluding, we emphasize a few key takeaways. First, LLM-powered ETL automation offers a significant leap in productivity for data engineering tasks, akin to how compilers and IDEs boosted software development in earlier eras. It allows engineers to work at a higher level of abstraction, focusing on what needs to be done rather than how to explicitly code it. Second, domain-agnosticism is a realistic goal thanks to foundation models' broad knowledge; an LLM can help with a marketing data pipeline one day and a genomic data pipeline the next, by virtue of its training on diverse information. This cross-domain capability is something traditional bespoke solutions lack. Third, the technology is ready to be applied in real-world systems, but success depends on integrating it thoughtfully – providing context to models, maintaining human checks, and logging and learning from the outcomes. Early adopters in industry (like the ones we cited) have shown it can work in production settings. We anticipate that in the near future, AI-assisted ETL will become a standard feature of data platforms, much like how AI code assistants are becoming common for developers.

Finally, we advocate for continued scholarly attention in this intersection of AI and data engineering. As foundation models evolve (becoming more accurate, acquiring larger context windows, or real-time learning abilities), their utility in pipeline automation will only grow. Future research should address open questions around verification of AI-generated pipelines, collaborative human-AI development interfaces, and how to scale these solutions for the largest of data environments. By addressing these, we can fully unlock the potential of foundation models to not just speed up ETL development, but to enable entirely new levels of adaptability and intelligence in data pipelines. In doing so, we move towards a future where organizations can integrate and utilize data faster, more reliably, and with greater insight than ever before – a true testament to the power of combining foundational AI with the foundations of data engineering.

REFERENCES

- [1] Bommasani, R. et al. (2021). On the Opportunities and Risks of Foundation Models. arXiv preprint arXiv:2108.07258 ([2108.07258] On the Opportunities and Risks of Foundation Models). (Introduced the concept of "foundation models," describing their broad capabilities and emergent properties.)
- [2] Santos, A. et al. (2025). Interactive Data Harmonization with LLM Agents. arXiv preprint arXiv:2502.07132 ([2502.07132] Interactive Data Harmonization with LLM Agents) ([Literature Review] Interactive Data Harmonization with LLM Agents). (Proposed the Harmonia system where an LLM-driven agent assists in building data integration pipelines for schema and value harmonization.)
- [3] Colombo, A., Bernasconi, A., & Ceri, S. (2025). An LLM-assisted ETL pipeline to build a high-quality knowledge graph of the Italian legislation. Information Processing & Management, Special Issue on Knowledge Graphs and LLMs (An LLM-assisted ETL pipeline to build a high-quality knowledge ...). (Demonstrated improved data quality in knowledge graph construction by integrating LLMs into the ETL process.)
- [4] Benson, G. et al. (2023). How SnapLogic built a text-to-pipeline application with Amazon Bedrock to translate business intent into action. AWS Machine Learning Blog (How SnapLogic built a text-to-pipeline application with Amazon Bedrock to translate business intent into action | AWS Machine Learning Blog) (How SnapLogic built a text-to-pipeline application with Amazon Bedrock to translate business intent into action | AWS Machine Learning Blog). (Describes SnapLogic's SnapGPT, which uses an LLM via Amazon Bedrock to generate integration pipelines from natural language, significantly reducing development time.)
- [5] Huthmacher, F. (2023). Supercharge Your ETL Pipeline: Fixing Data Quality Issues with AI. Medium (Supercharge Your ETL Pipeline: Fixing Data Quality Issues with AI | by Felix Huthmacher | Medium) (Supercharge Your ETL Pipeline: Fixing Data Quality Issues with AI | by Felix Huthmacher | Medium). (Explains an approach to integrate a Bedrock-hosted LLM (Claude) into AWS Glue to automatically correct data quality issues during ETL, including architecture and code samples.)
- [6] Informatica. (2023). Building ETL Pipelines with AI (Building ETL Pipelines with AI | Informatica) (Building ETL Pipelines with AI | Informatica). (Discusses how generative AI can automate parts of ETL pipeline development, including code generation, recommendations, and low-code solutions, from a data management industry perspective.)
- [7] Pandey, B. K. (2025). AI-Powered ETL Pipeline Orchestration: Multi-Agent Systems in the Era of Generative AI. ODSC Webinar/Medium Article (AI-Powered ETL Pipeline Orchestration: Multi-Agent Systems in the Era of Generative AI | by ODSC - Open Data Science | Feb, 2025 | Medium) (AI-Powered ETL Pipeline Orchestration: Multi-Agent Systems in the Era of Generative AI | by ODSC - Open Data Science | Feb, 2025 | Medium).





[www.ijirset.com |A Monthly, Peer Reviewed & Refereed Journal| e-ISSN: 2319-8753| p-ISSN: 2347-6710|

Volume 14, Issue 4, April 2025

|DOI: 10.15680/IJIRSET.2025.1404004|

(Introduces the concept of agentic ETL pipelines using multiple LLM-based agents for tasks like data retrieval, transformation, and loading, enabling self-healing and adaptive workflows.)

- [8] AWS. (2023). Amazon Bedrock Developer Guide (Supercharge Your ETL Pipeline: Fixing Data Quality Issues with AI | by Felix Huthmacher | Medium). (AWS documentation describing Amazon Bedrock, a managed service for accessing foundation models via API, used in our implementation to ensure secure and scalable LLM integration.)
- [9] O'Brien, C. et al. (2023). LLMOps: Operationalizing Large Language Models. Databricks Blog (LLMOps: Operationalizing Large Language Models Databricks). (Covers best practices for deploying and monitoring LLMs in production, relevant to managing the LLM components of ETL pipelines for reliability and performance.)
- [10] Amazon CodeWhisperer Announcements. (2022). (How SnapLogic built a text-to-pipeline application with Amazon Bedrock to translate business intent into action | AWS Machine Learning Blog) (Illustrative reference to foundation model applications in code generation, showing that models can translate natural language to code in languages like Python/SQL, analogous to pipeline generation tasks.









INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY





www.ijirset.com