

A Survey on Scalable Feature Engineering Techniques for Cloud-Native Machine Learning Workflows

Raviteja Narra

Independent Researcher

ravitejanarra563@gmail.com

Abstract: *The presented survey provides a briefing of scalable feature engineering techniques of cloud-native machine learning (ML) processes based on the findings of the recent literature. It explores distributed feature extraction, large-scale feature transformations, feature selection, dimensionality reduction and feature representation using deep learning. The automated practices include stores, real-time streaming, feature versioning, feature lineage and feature governance are explored with respect to their perceived effects on the activity of reproducibility, consistency and scalability. The literature presents high-dimensional skewed data, computational constraints, orchestration complexity, and low-latency as some of the key Challenges in the field. The survey gives an overview of commerce offs between throughput and latency and resource efficiency and discusses integration approaches to scattered processing systems using cloud-native achievement systems. The observations can be utilized during the designing of adaptive feature engineering pipeline, resistant, and efficient within the modern cloud environment.*

Keywords: Scalable Feature Engineering, Cloud-Native Machine Learning, Automated Feature Engineering, Real-Time Feature Engineering, Resource Management, Dimensionality Reduction, Deep Learning Feature Representation

I. INTRODUCTION

The widespread application of cloud-native architectures has entirely transformed the architecture and the functioning of contemporary machine learning (ML) systems. Cloud-native systems consist of distributed, scalable, highly dynamic tools that consist of microservices, containers and orchestration systems that support scalable and resilient ML processes [1]. In such settings, feature engineering has emerged as a significant factor in the model performance, computation efficiency and scalability of pipelines, although large and heterogeneous data are handled with special issues, without compromising consistency and efficiency [2].

Feature engineering includes cleaning, normalization, construction, and the selection of pertinent features from raw data as a way of supporting learning [3][4]. With the increasing volume of data and dimensionality in fields like text mining, genomics, sensor analytics, and others, scalable feature engineering algorithms are becoming more and more needed to combat redundancy, noise and computational bottlenecks [5]. In comparison to the low-dimensional datasets, more recent data-driven systems are being run on high-dimensional and heterogeneous data, and this is aggravating the relevant feature problem, feature redundancy, and feature overfitting [6].

The most important feature engineering task is feature selection, which is intended to find the most informative set of features. Due to its NP-hardness and combinatorial complexity, feature selection is commonly solved by means of metaheuristic or evolutionary algorithms that can give scalable approximations in large scales [7][8]. Besides feature selection, novel forms of advanced automation like NAS, HPO and automated model selection have become increasingly popular, allowing scalable and adaptive ML processes and deemphasizing the need to resort to manual design selections. The feature engineering pipelines in cloud-native environment should be co-located with the distributed data processing

systems, workflow orchestrators and auto-scalers to offer batch and streaming workloads. Scalability, consistency and latency of training-serving are also important concerns and are of special concern to dynamic and multi-tenant cloud environments [9]. Current survey articles are focused on the distributed ML infrastructure, MLOps, or feature engineering optimisation, but not the entire image of scalable feature engineering, in particular, and practicable tools, architectural principles, and policies of automation. [10].

The presented paper fills this gap with a systematic review of scalable feature engineering techniques to ensure cloud-native ML workflows are introduced. It explores the distributed processing paradigms, feature transformation and selection methods, automation and integration with feature stores strategies, bottlenecks, trade-offs, and gaps in research. The goal is to help practitioners and researchers construct efficient, scalable, and production-ready ML pipelines that are able to support modern and high-dimensional datasets.

A. Structure of the Paper

The structured of the paper is as follows. Section II introduces cloud-native machine learning workflows, emphasising distributed data processing frameworks, data pre-processing, and resource management strategies. Section III addresses scalable feature engineering techniques and associated challenges in large-scale environments. Section IV discusses automated and intelligent feature engineering, focusing on feature stores, real-time processing, and governance mechanisms. Sections V and VI provide a summary of the literature review, followed by conclusions and directions for future work.

II. CLOUD-NATIVE MACHINE LEARNING WORKFLOWS

Cloud native machine learning workflows use distributed data processing systems, scale control systems and auto-scaling systems to efficiently construct, train and deploy AI models on a large scale. Apache Spark and Apache Flink frameworks make it easy to provide high-throughput processes in batch and stream processing and form the foundation of data-intensive ML pipelines. Pre-processing and feature engineer services use cloud-based systems and serverless and distributed computing models to make the process of scale pipeline organization easier, and achieve scalability, and fault tolerance. Clouds also allow for dynamically changed ML system performance to workload needs, cost-effective performance, and resource utilisation through the use of auto-scaling, resource allocation, and advanced load-balancing algorithms.

A. Data Processing Frameworks for Cloud Environments

Hybrid processing frameworks are versatile and can handle both batch and stream processing, among other data types. In the cloud-native ML pipelines, models are mainly trained using more frameworks like Apache Spark, Apache Flink and Apache Beam, instead of training models to perform feature extraction, transformation and aggregation steps:

1) Apache Flink

Apache Flink is an up-to-date platform for streaming analytics and distributed processing. It is a next-gen framework for processing massive amounts of data, and it works with all typical cluster configurations with low latency and great throughput. Its applicability to scalable feature engineering is that it can compute features in real-time that are based on unlimited data streams with event-time semantics and stateful processing. Flink currently supports unlimited and limited data processing using its Batch and DataStream APIs, and allows one to build hybrid feature pipelines that integrate both historical and real-time features. Although resource management is not built-in in Flink, it can be integrated with cloud-native orchestration systems to achieve real-time scaling of streaming feature pipelines, so it can be used in time-sensitive applications like monitoring and anomaly detection [11].

2) Apache Spark

Apache Spark is a free and open-source platform designed to handle massive data sets efficiently and for advanced analytics. AMPLab at UC Berkeley developed it in 2009, and it's user-friendly. Spark makes it simple to build applications using Python, Scala, or Java. In feature engineering pipelines Spark is mostly used to compute features in

batch tasks which include large-scale aggregation, joins, and feature normalization. Its ability to support distributed Data Frame operations and in-memory execution can be useful in the generation of high-dimensional feature sets using large datasets. Spark ML and MLib have pipeline abstraction, which enables the integration of feature preprocessing with model training to be used downstream. In-memory computation model of Spark is much more scalable and efficient compared to disk-based processing frameworks in offline feature engineering pipelines [12]. Spark is usually used as a core feature processing engine in cloud-native ML systems owing to these features.

3) Apache Beam

The open-source programming model known as Apache Beam is mostly used for data processing and is designed to be both scalable and portable. Its main contribution to scalable feature engineering is that it is possible to specify feature extraction and pre-processing logic once and run it in many distributed execution engines, including Apache Spark and Apache Flink [13]. Beam can support more than just windowing, event-time processing, and stateful computations that are essential in creating a consistent feature engineering pipeline both in batch and streaming systems. Consequently, Apache Beam can be easily used with cloud-native ML workflows, which need scalable, versatile, and cross-platform feature engineering.

Table I gives a comparative description of Apache Spark, Apache Flink, and Apache Beam with reference to their ability to sustain feature engineering workloads in a scalable manner

Table 1: Comparison of Data Processing Frameworks for Scalable Feature Engineering

Framework	Processing Model	Latency Characteristics	State Handling	Feature Engineering Suitability
Apache Spark	Primarily batch, micro-batch streaming	Medium to high (batch-oriented)	Limited native state handling (batch-centric)	Well-suited for offline feature engineering, large-scale aggregations, joins, and historical feature computation
Apache Flink	Native streaming with batch support	Low latency (event-driven)	Strong stateful processing with event-time semantics	Ideal for real-time and online feature engineering, continuous feature updates, and time-sensitive applications
Apache Beam	Unified batch and streaming abstraction	Depends on execution engine	Abstracted state and windowing model	Suitable for portable and unified feature pipelines across batch and streaming environments

B. Data Preprocessing and Feature Engineering in the Cloud.

The raw data must be processed and then converted into a suitable format before it can be used to train AI models. Cloud-based AI systems provide several data preprocessing and feature engineering tools and services, which include:

1) Data Transformation Services

Services offered by cloud providers such as AWS Glue, Azure Data Factory, and Google Cloud Dataflow enable customers to construct, organise, and manage data pipelines for the processing and improvement of large data sets. These managed services will minimize overhead in the operations and offer scalability, monitoring, and fault tolerance, yet could restrict the ability to perform custom transformations.

2) Serverless Computing

Users can run data preparation and feature engineering jobs without managing the underlying infrastructure using serverless services like AWS Lambda, Azure Functions, and Google Cloud Functions [14]. These systems can handle massive volumes of data efficiently and at a low cost by autonomously scaling up or down depending on the workload. Nonetheless, serverless can cause latency when performing large and complex computations of features and should be coordinated with caution to prevent inconsistencies between parallel computations.

3) Distributed Data Processing Frameworks

Distributed data processing is a feature of Apache Spark and Apache Flink that allows users to preprocess and convert

large amounts of data efficiently using cloud resources. Microsoft HDInsight, Amazon Elastic Medical Record, and Google Cloud Datapost are managed services that can incorporate such frameworks. While open-source frameworks provide full control and flexibility for custom feature engineering, they demand careful resource and job management. Decisions made during the preprocessing—such as normalization, categorical encoding, or feature selection—must be consistent across training and serving to prevent feature skew that can reduce the model performance.

C. Load Balancing and Resource Allocation Strategies

Efficient load balancing and resource allocation can be used to optimize the AI systems running on clouds especially where scalable feature engineering pipelines are to be run with large amounts of data to be pre-processed and transformed efficiently. The major approaches and influencing aspects include the following:

1) Horizontal Scaling

This policy is associated with the addition or removal of computing assets, like virtual computers or containers, to scale back the workload and maximise the utilisation of resources [15]. In case of feature engineering pipelines, horizontal scaling makes it possible to execute feature processing functions, including feature extraction, encoding, and aggregation, in parallel and thereby minimizes end-to-end latency and maximizes throughput. However, scaling out introduces cost and coordination overheads, especially in cases where skewed distributions of features cause straggler tasks thus there is a need to weigh performance benefits with the cost of operation. Auto Scaling Groups on AWS, Virtual Machine Scale Sets on Azure, and Cloud Instance Groups on Google are all examples of provider-native capabilities that can be utilised for horizontal scaling.

2) Vertical Scaling

Increasing or decreasing the processing power of specific components, such as a central processing unit, RAM, or graphics processing unit [16]. Vertical scaling can accelerate compute-intensive feature engineering tasks, such as large matrix transformations, high-dimensional feature construction and deep feature generation. The limitation of this strategy lies in the hardware of individual instances limit its usefulness in very large-scale or very parallel feature pipelines.

3) Resource Allocation Policies

Resource allocation policies, including the CPU and memory quotas, may be used to make sure that the resources are evenly distributed among several AI workloads and to avoid resource contention. In feature engineering pipelines, carefully designed allocation policies ensure that resource-intensive pre-processing stages do not starve downstream tasks such as model training or inference, thereby maintaining pipeline throughput and controlling cloud costs. Cloud service providers facilitate the specification and implementation of such policies by means of such tools as AWS Resource Groups, Azure Resource Manager, and Google Cloud Resource Manager.

D. Auto-scaling and Dynamic Resource Provisioning

Auto-scaling and dynamic resource provisioning enable cloud-based AI systems to automatically scale the resources allocated to them based on the demands of the workload, thereby enhancing efficiency, decreasing expenses, and having a direct effect on the performance of feature engineering pipelines by controlling latency and throughput. Key methods include:

1) Reactive Auto-scaling

Resource allocation in reactive auto-scaling occurs in reaction to predetermined triggers, while metrics system indications such as CPU load, memory consumption, or request response time are monitored. In feature engineering pipelines, reactive auto-scaling ensures that data pre-processing workload spikes (e.g. abrupt increases in batch size or feature transformation) do not form bottlenecks, which sustain steady throughput and eliminate latency spikes.

Some of the services offered by cloud suppliers are AWS Auto Scaling, Azure Auto Scale, and Google Cloud Autoscaler, which allow users to implement reactive auto-scaling policies.

2) Predictive Auto-scaling

Predictive auto-scaling is a method that uses machine learning to forecast future workload demands by analysing

historical workload trends. In large-scale feature extraction tasks, predictive scaling enables systems to allocate resources to meet peak loads, which reduces processing delays and provides high throughput on tasks to compute features in one-hot encoding or feature normalization. Google Cloud AI Platform, Azure ML, and AWS Forecast can all be used to implement predictive auto-scaling.

3) Serverless Computing

Serverless computing platforms (such as AWS Lambda, Azure Functions, or Google Cloud Functions) automate resource scaling in response to workload demand, doing away with the need for manual provisioning. The serverless architecture specifically can be useful in event-driven feature engineering workloads, including real-time data transformations or feature aggregation, and is being used by auto-scaling to keep the latency low even under varying load..

III. SCALABLE FEATURE ENGINEERING TECHNIQUES

Advanced feature engineering entails deriving useful representations of raw data with domain knowledge, statistical analysis or machine learning. In cloud-native environments, the focus shifts from technique definitions to scalability, performance bottlenecks, and distributed execution. Domain-specific features represent problem-specific signals, whereas statistical methods, such as PCA, factor analysis and clustering, indicate latent structure. Time-series features (lag variables, rolling statistics, frequency-domain transforms) are used to capture temporal dependencies and embeddings or text, image, and signal-specific extraction functions are used to represent the data. On large scale, feature engineering uses parallel, vectorized, and GPU-accelerated computations, which makes trade-offs in the cost of computation, latency and predictive performance in an iterative search.

A. Distributed Feature Extraction.

Distributed feature extraction involves sharing the workload across numerous nodes in a distributed computing environment in order to compute and generate features from a large-scale dataset. Some examples are group-by aggregations (e.g. mean sales per store), rolling/windowing features (e.g. 7-day moving averages), and one-hot encoding at high cardinality. The operations are implemented concurrently through frameworks such as Apache Spark and Flink. Scaling challenges include skewed data distributions causing stragglers, memory pressure from high-cardinality features, and maintaining consistency across streaming and batch jobs.[17][18]. Distributed feature extraction is also more fault-tolerant and scalable; feature pipelines can dynamically adjust to varying data loads and remain consistent and reproducible across both batch and streaming workloads.

B. Feature Transformation at Scale.

Categorical features are variables that can only have a certain number of names, or groups. Ordinal data and nominal data are two types of categorical data that are derived from nature. When there is an inherent hierarchy in the values of an attribute, say that it is ordinal. Representing categorical data in numerical form at scale presents some challenges that include cardinality explosion and memory pressure as well as network overhead in distributed systems. Although conventional encodings are popular, the issue of scalability has to be taken into consideration.

1) Label Encoding

Label encoding transforms labels that do not contain numbers into numerical values. It can be used on ordinal features although it can add unwanted order biases. Categorical variables can take significant memory during replication in distributed pipelines with large numbers of them. Current substitutes such as target encoding are able to minimize the memory footprint without compromising predictive relationships [19]. It is still possible to perform label encoding with Label Encoder through Python scikit-learn.

2) One Hot Encoding

The process of one-hot encoding involves creating dummy columns for every possible category of an attribute. This may result in a memory explosion and more shuffling operations in a distributed system such as Spark in high cardinality, which affects throughput and latency. The problems can be addressed by hash features or sparse representations that

enable parallel transformation in large scale.

3) Binary Encoding

The binary process of encoding starts by converting the values of classes numerically and then the textual representation of the values. The second one is splitting the binary numbers into their columns as shown in Figure 1. This reduces the size compared to one-hot encoding and reduces the memory load of a distributed system.

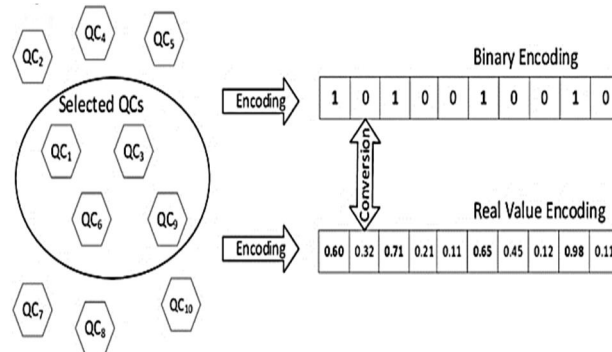


Fig. 1. Binary Encoding

4) Modern Feature Transformation Techniques

In feature engineering pipelines on a large scale, techniques like target encoding, feature hashing and quantile-based encoding are gaining popularity as well. The methods will make sure the per-dimensionality of categorical variables is reduced, the network communication is reduced, and scalable transformations of cloud-native machine learning processes may be delivered. Elevated stage feature representations are also significant methods in streaming and batch processing pipelines to ensure latency is minimal, and throughput is elevated.

C. Feature Selection and Dimensionality Reduction

The aim of feature selection is to minimise computation resources whilst maximising prediction performance by finding the minimum number of features to do so. The procedure manages to reduce complexity of the final model by eliminating the unnecessary features. Useless features tend to decrease performance, hence overall performance can be significantly enhanced [20]. The feature selection also helps determine what set of features can best be combined with the ML model and may shed some light on other important performance determinants. Furthermore, a simplified model can be created by removing features to the most meaningful ones. Therefore, feature selection has four significant roles:

- Reduction of the size of a model by deleting redundant and unhelpful features, resulting in fewer computations and less memory requirements.
- The probable enhancement of performance due to the removal of potentially counterproductive features.
- Reduced feature set leads to simplification of the model structure.
- Knowledge discovery (i.e. gives information about the inner representation)
- Data interpretability and interpretation of the resulting machine learning model [21].

In the distributed cloud environment, feature selection should consider the computational overhead, data divisions, and communication expenses. Filter-based (e.g., mutual information) and embedded (e.g., Lasso, tree-based) methods are better than wrapper and recursive methods. Dimensionality reduction PCA (SVD, autoencoders) simplifies and reduces features and keeps the important information intact. Distributed implementations (e.g., Spark MLlib PCA, autoencoders) minimize data movement and memory bottlenecks, whereas standard SVD may need approximations at scale. These methods reduce feature space, enhance training, reduce the cost of computation and storage and overcome overfitting in large and heterogeneous data.

D. Feature Construction and Generation

Feature construction at scale entails the automatic generation and assembling of new features out of raw and intermediate data whilst maintaining an efficiency, reproducibility and reliability over distributed pipelines. Manual or hand-written extractors of features are not feasible in cloud-native applications because of large volumes of data, non-homogeneous data types, constant schema changes and the high frequency of iterative experiments. AutoML-based feature generation systems and feature synthesis libraries are also rapidly gaining popularity in accessible frameworks of automated generation of feature constructions, so as to perform systematic search of feature combinations and transformations with minimal human intervention. Such systems provide rapid development speeds at the tradeoff of fresh burdens of explosion of search space, scale-based computational overhead and redundancy of features.

In distributed environments, feature construction should be well coordinated on many nodes to prevent stragglers and excessive data exchange, and an uneven feature cardinality and skew of data can cause a significant effect on load balance and pipeline performance. Moreover, it also becomes important to control the intermediate feature materialization and storage, the uncontrolled materialization of features might consume more memory and I/O resources, and eventually result in the impact on throughput and end-to-end latency in large-scale feature engineering pipelines [22].

E. Deep Learning-Based Feature Representation.

Deep learning algorithms consist of neural networks that can learn to extract useful low-level data from raw data, with little human supervision. Similar to the biological brains, deep learning pipelines process data based on huge networks of interconnected layers to extract semantic features on many levels [23]. An example of this is computer vision in which shallow layers are used to absorb low-level visual concepts such as shapes, lines, and edges, and deep layers absorb higher-level semantics [24]. The feature extraction is performed by means of deep learning, which adds further concerns to distributed, cloud-native settings. Big data relocation among nodes and between nodes and GPUs can increase training expenses whereas memory and data bandwidth bottlenecks are also able to influence throughput. The representation of deep features in production should be orchestrated to ensure that model inference is carried out across distributed systems to reduce the latency. Such methods as model parallelism, data parallelism, and storing intermediate feature maps alleviate them and guarantee scalable and efficient feature representation in cloud-native ML workflows.

F. Data Preprocessing Requirements for Deep Learning Models

Deep learning models rely significantly on high-quality training data for their performance due to the high level of automation in their data processing. But in many real-world cases, the needs of the intended application or model do not align with the qualitative qualities of the raw data. To that end, data preprocessing is now a crucial step in creating deep learning applications [25]. Basic data preparation tasks that comprise preprocessing include cleaning (the removal of outliers and inconsistent data points), encoding categorical values, imputation of missing values, and labelling.

1) The Need for Data Augmentation

The acquired data might not be adequate or representative enough for the intended purpose in many Big Data applications. Data augmentation is employed in such cases after initial preprocessing to generate additional samples by modifying the current data, thus increasing the quantity and variety of the training data. When it comes to improving deep learning models' ability to generalize, this is by far the most common method. Issues including bias, unbalanced data, and domain shift have also been tackled by data augmentation.

2) The Need for Feature Engineering In ML

ML models trained on DNNs can pick up helpful features on their own, but explicit feature processing is still necessary for certain tasks to get good results. Applications that require sophisticated data and are scarce, such as industrial applications, credit risk evaluation, online fraud detection, and machinery problem diagnosis, fall under this category. There is evidence that deep learning networks with dedicated feature extractors perform much better by producing more accurate feature representations of the input data. The data is typically massive, diverse, complicated, and noisy, necessitating the management of massive feature sets, frequently containing duplicate and redundant samples. The three

main steps of machine learning—data pre-processing, data augmentation, and feature engineering—are shown in Figure 2.

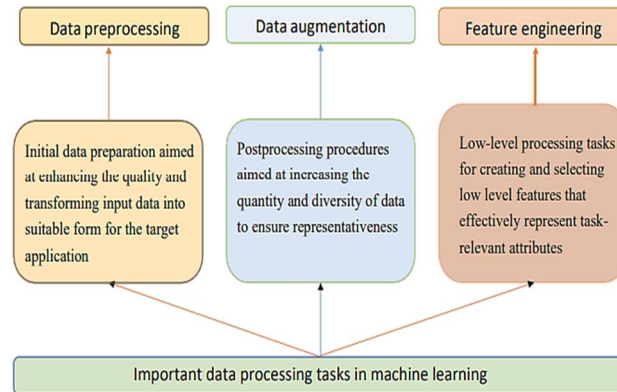


Fig. 2. Data Processing Tasks in Machine Learning in Cloud-Native

These steps transform raw data into representations that the model can learn from. The mentioned stages not only enhance the quality of the data but also increase its diversity and the extraction of attributes relevant to the task.

G. Scalability Challenges in Feature Engineering

The automated data processing method also has a serious problem with scalability. These methods excel at solving large-scale problems with massive datasets, but they bomb when faced with small data. Training, on the other hand, is a breeze with big datasets. When data is limited, manual processing approaches usually yield better results than their automated counterparts. In addition, there is a direct correlation between the computational expense of automated processing methods and the size of the dataset, the nature of the problem, and the resources available for processing:

1) Data Volume and High Dimensionality

Among the main scalability issues is the fact that the amount of information produced by the distributed systems, IoT devices, and online platforms is exponentially growing [26]. High-dimensional feature spaces can frequently result in large datasets, which make their extraction, transformation, and selection more computationally complex and more memory-intensive. The storage and parallel computation of such data are best handled by distributing it, which makes pipeline design and implementation more difficult.

2) Computational and Resource Constraints

Scalability of feature engineering requires substantial computational resources, particularly when more complex transformations, such as feature interactions, feature embeddings, and dimensionality reduction, are needed. Unproductive use of resources may result in more latency and operational expenses in clouds. The issue of striking a balance between performance and cost-efficiency and dynamically allocating the CPU, memory, and graphics cards is a significant problem.

3) Data Skew and Distributed Processing.

Feature engineering in cloud-native systems is typically implemented on distributed nodes. Data skew (or uneven distribution of data) may lead to imbalance in the workload resulting in straggled data and poor system performance. To guarantee scalability and reliability, it is necessary to ensure that the data partitioning is evenly distributed and the processing tasks are scheduled effectively tributed Processing and Data Skew.

4) Feature Consistency and Reproducibility

It is not easy to be consistent in the features employed when training models and the ones used in inference at scale. Mismatch between the source of data, logic of transformation or the environment where the data is executed may lead to training serving skew which negatively impacts the model accuracy. It is thus not nontrivial to version, keep lineage records and guarantee reproducibility of feature pipelines in large-scale deployments.

5) Real-Time and Low-Latency Requirements

A lot of contemporary applications, such as recommendation systems and fraud detection, require a lot of real-time or near-real-time feature computation. The inability of feature engineering pipelines to scale to satisfy very strict latency guarantees but handle high-velocity data streams adds further complexity and, in particular, integrating batch and streaming workflows.

6) Workflow Orchestration and Fault Tolerance

Scalable feature engineering pipelines can be made of many interdependent tasks that coordinate between distributed infrastructures. Failure can be caused by crashing of nodes, network problems or contention on resources which can interrupt feature generation processes. Fault-tolerant workflow design and recovery mechanisms are both fundamental but difficult in dynamically operating clouds.

7) Governance, Security, and Compliance

When feature engineering pipelines grow over organizational borders and cloud software, the process of ensuring secure access to data, maintaining privacy, and adhering to regulations becomes more complicated. The ability to apply scalable governance without sacrificing performance is an open challenge. Feature space complexity, dimensionality reduction is significantly more efficient in model training, reducing the cost of computation and storage, and addressing overfitting, so is particularly useful on large, heterogeneous datasets trained on distributed clouds.

IV. AUTOMATED AND INTELLIGENT FEATURE ENGINEERING IN CLOUD-NATIVE

Automated and intelligent feature engineering are two practices that aim to handle everything automatically and make feature pipelines more cloud-native, with better scalability, consistency, and adaptability. An important benefit of automation is that it enables a consistent and methodical process of feature development and selection across the whole dataset, regardless of its size or heterogeneity. The use of cloud-native enables the sharing of features, instantaneous computation of features, and the integration of training and testing without disturbance, through feature stores, stream processing systems, and metadata-driven orchestration.

Intelligent and automated feature engineering seeks to reduce the human factor and enhance the scalability, uniformity, and flexibility of feature pipelines through cloud-native machine learning workflows. These solutions also facilitate immaculate reuse, real-time computation of environmental features and synchronisation of training and inference using feature stores, stream processing, and metadata-based automation. Smart feature management also brings with it versioning, lineage tracking and governance so as to provide reproducibility, compliance and operational reliability. All these developments are a platform for healthy MLOps and scale continuous learning systems.

A. Feature Stores and Feature Management Systems

Feature stores provide a single point of storage to store and retrieve features and ensure consistency between training and serving machine learning to aid in preventing training-serving skew. Conventional data workflow models were highly model and environment-specific by nature which in effect resulted in data drift, redundancy and inefficiencies. The feature stores are a reversal of this paradigm, both in that feature logic and model logic are decoupled, and that there is an API to access the same features during training (offline) and during inference (online). Real-time ML pipelines need feature stores to take in event sources, like Kafka, Kinesis, or Pulsar, and transform them into useful features, and offer them with low latency [26]. This immediacy is needed in such applications as fraud detection, real-time recommendations, anomaly detection, and autonomous systems. The complexity is added by the dual-read / write scheme of access which supports both the working loads of the batch and the streaming. Offline stores are typically placed in scalable warehouses such as BigQuery, Redshift or data lakes in contrast to online stores that are typically fast key-value stores such as Redis, Cassandra or DynamoDB. Popular feature store implementations have varied architecture and application: e.g., Feast is flexible, cloud-neutral, Tecton is all real-time feature computation and monitoring, whereas Hopsworks combines feature lineage and governance. Such variations help practitioners to choose those stores that are tailored to some training and service requirements. Futuristic coordination of these two stores guarantees consistency

and freshness that are key in high-performing models.

B. Streaming and Real-Time Feature Engineering.

Stream pipelines need to have a frequent connection between streaming processing systems and the feature store. Stream events of operational databases, telemetry systems or user activities are consumed into stream processors, where they are transformed and then deposited in the feature store [27]. The integration of this is not easy particularly when one is dealing with state management, backpressure and out-of-order event handling. Streaming state, which is vital to correct feature computation, is operated by frameworks such as Flink, and is provided with windowing and checkpointing functionality, as well as the absolutely indispensable exactly-once semantics. After computations on features are made, they should be written to low-latency online stores. There are systems that use Change Data Capture (CDC) patterns to transfer updates of stream processors orthogonally to online stores in near real time. Message queues are used by others to decouple storage and compute.

A key challenge is coordination of offline and online stores. As an example, consider a real-time recommendation model in which a feature is defined as user clicks within the last hour. The online store is constantly updated with streams, but the offline store (to be used in training) could be delayed. To ensure consistency, the scheduled batch jobs recreate the recent streaming information into the offline store such that they match the online calculations to ensure that the model is trained on the correct historical characteristics. The feature values should be time-stamped to enable point-in-time joins, particularly when training a model with a historical background. Such metrics of data freshness as staleness thresholds and freshness SLAs are essential.

C. Feature Versioning, Lineage, and Governance.

A scalable feature store is not only a store of values it is a store of lifecycle of features. This involves defining, registering, modifying, depreciating as well as deleting features whilst ensuring traceability and reproducibility. The versioning of features is required when the definition of features is evolving or when models are built on top of other versions of logic. It possesses a good versioning that annotates feature changes with semantic metadata, dates and author information. It allows the reversal of the past definitions and provides backwards compatibility of the models when manufacturing. Extraction of lineage provides a query able map in a pictorial display of how the features of the data on the source were transformed into the final products. It is also used to determine the effects of changes in data schema and is also necessary to guarantee regulatory traceability, which is especially relevant to compliance regimes, including GDPR, HIPAA, and SOC [28]. The policies of data governance have to be documented in the feature registry. As an example, the labelling of a feature as PII can cause masking policies, encrypting requirements, or access controls. Similarly, labelling features as financial- or healthcare-sensitive triggers auditing and retention policies. Documentation and discoverability are also critical.

V. LITERATURE REVIEW

This section is a review of previously conducted work, which backs scalable feature engineering in cloud-native machine learning workflows, either directly or as facilitated by supporting systems, like workflow orchestration, distributed execution, and MLOps infrastructure. The literature is categorized into four, namely workflow orchestration, hybrid cloud HPC pipelines, scaling benchmarking, and cloud-native ML automation, each of which has considered major scales of scalability issues in large-scale feature pipelines.

The cloud-native workflow orchestration researches present fundamental support of scalable feature engineering through controlling the multi-stage, feature-extraction, and transformation activities. Shan et al. (2023) present Kub Adaptor, a cloud-native workflow engine that is intended to connect workflow management systems to Kubernetes. It applies to scalable feature engineering since feature pipelines can be composed of processing stages that are interdependent, so proper and repeatable execution of feature pipelines requires a regular schedule and fault-tolerant orchestration [29].

Olaya et al. (2023) introduce a way of running a scalable scientific workflow in hybrid cloud and HPC systems that is

illustrated with a workflow of soil moisture prediction with the use of ML. Their system facilitates scalable feature engineering by spreading the huge feature pre-processing and data preparation tasks over heterogeneous computing resources [30].

Shin et al. (2022) focus on workflow-conscious scheduling to optimize the use of resources and to speed up the time to complete a workflow. These types of scheduling are also essential to scalable feature engineering pipelines, in which costly steps like aggregations and joins can be affected by stragglers and contention of resources in a distributed system [31].

Henning and Hasselbring (2022) present a standard benchmarking approach to assessing the scalability of cloud-native applications. The work could be used to evaluate scalability of feature engineering pipelines by giving metrics to measure the throughput, latency, and resource efficiency at increasing data and workload scales [32].

Peña-Monferrer, Manson-Sawko and Elisseev (2021) offer a cloud-based hybrid system that analyses and simulates drop dispersions efficiently. A cloud-native framework processes CFD outputs, while a HPC cluster runs the simulation component. The study is based on CFD simulations, but its pipeline architecture with modules and the principles of flexible resource management can be directly applied to scalable workflows of feature engineering[33].

George et al. (2020) present Katib, a cloud-native Kubernetes-based system that serves to optimize hyperparameters. Though Katib mostly focuses on model optimization, it can be used to augment scalable feature engineering because it facilitates automated experimentation and feature pipelines on cloud-native MLOps systems [34].

Table II provides a summary of the reviewed papers and their relevance to scalability and infrastructure support to cloud-native feature engineering processes.

Table 2: Scalable Feature Engineering Techniques for Cloud-Native Machine Learning

Authors	Focus Area	Objectives	Approach	Key Findings	Future Work
Shan et al. (2023)	Cloud-native workflow orchestration with Kubernetes	Kubernetes task scheduling and workflow scheduling algorithms must be compatible with one another	A universal docking framework called KubeAdaptor was proposed to integrate workflow systems with Kubernetes. The architecture design included features like event-driven execution, fault tolerance, and containerisation.	Utilised four operational scientific workflows to prove efficacy; established consistent order for task scheduling across Kubernetes and workflow engines.	Extending support to heterogeneous schedulers, improving adaptive scheduling, and optimizing performance for large-scale multi-tenant environments
Olaya et al. (2023)	Scientific workflows on cloud and HPC platforms	To enable scalable scientific workflow composition and execution across cloud-native and HPC environments	Designed a workflow orchestration approach integrating HPC (LSF), Kubernetes (K8s), and object storage; validated using the SOMOSPIE ML-based soil moisture prediction workflow	Achieved scalable deployment of ML models across a range of spatial resolutions and sizes; shown adaptability across hybrid infrastructures	Automation of workflow optimization, improved data locality strategies, and broader application to other scientific domains
Shin et al. (2022)	Workflow-aware task scheduling	With the goal of maximising efficiency and reducing the time it takes to	Created a scheduling algorithm that takes workflows into account by combining convex task-splitting techniques	Better balanced workflow completion time with resource utilisation; hence,	Integration with real-world workflow engines, dynamic workload adaptation, and

		complete weighted workflows,	with hybrid priority rules; tested it in a simulated environment.	outperformed deterministic parallelism.	support for heterogeneous cloud resources
Henning & Hasselbring (2022)	Testing the scalability of apps built for the cloud	To provide a standardized method for empirical scalability evaluation	Introduced a benchmarking methodology including scalability metrics, measurement techniques, and a benchmarking tool architecture aligned with SLO-based evaluation	Enabled reproducible and statistically grounded scalability assessments for cloud-native systems across varying load and resource configurations	Development of automated benchmarking pipelines, integration with CI/CD, and energy-aware scalability metrics
Peña-Monferrer et al. (2021)	Hybrid HPC–cloud workflows for CFD simulations	To achieve agile and scalable CFD simulation analysis using hybrid infrastructures	Combined HPC-based CFD simulation with cloud-native microservice-based data analysis pipelines	Demonstrated dynamic scalability, reduced idle times, and efficient resource sharing across simulations and locations	Expansion to real-time analytics, tighter coupling between simulation and cloud services, and AI-driven optimization of CFD workflows
George et al. (2020)	Cloud-native hyperparameter tuning systems	To design a scalable, fault-tolerant, and multi-tenant hyperparameter tuning platform	Proposed Katib, a Kubernetes-native hyperparameter tuning system supporting multiple ML frameworks; evaluated through experiments and production deployments	Showed high scalability, extensibility, and robustness across on-premise and cloud environments; supported diverse user and administrator needs.	Enhanced AutoML integration, smarter search algorithms, and tighter coupling with MLOps pipelines

VI. CONCLUSION AND FUTURE WORK

Scalable feature engineering is an important enabler to cloud-native machine learning workflows that can be used to extract distributed features, perform transformations at scale, dimensionality reduction, and deep learning-based feature representations. The distributed processing structures, cloud-based preprocessing services and dynamic resource management are examined and suggest potential improved scalability, fault tolerance and performance. But there are a number of issues, such as skewed data, consistency of features between training and serving, low-latency demands, reproducibility, and control in the multi-tenant setup. Also, feature engineering in the form of automated and intelligent feature engineering, feature stores as building blocks, and with metadata-driven governance as its foundations, has a bright future, but it should be introduced carefully in production, bearing in mind trade-offs in resource usage, latency, and reliability of the service. Future research should be on adaptive feature engineering based on online learning, auto-optimal feature pipelines, explainable automated feature generation, privacy-preserving computation, less-energy feature

processing, and more intimate feature stores with foundation models. These open challenges need to be overcome with the help of strong, efficient, and practice cloud-native ML systems.

REFERENCES

- [1] P. Chandrashekar, "A Survey of Tools, Techniques, and Best Practices: CI/CD Integration in DevOps Workflows," *Int. J. Adv. Res. Sci. Commun. Technol.*, vol. 3, no. 3, pp. 1366–1376, Jul. 2023, doi: 10.48175/IJARSCT-11978V.
- [2] M. Mekki, N. Toumi, and A. Ksentini, "Microservices Configurations and the Impact on the Performance in Cloud Native Environments," in *2022 IEEE 47th Conference on Local Computer Networks (LCN)*, IEEE, Sep. 2022, pp. 239–244. doi: 10.1109/LCN53696.2022.9843385.
- [3] R. S and G. Battineni, "A Survey on Recent Advancements in Auto-Machine Learning with a Focus on Feature Engineering," *J. Comput. Cogn. Eng.*, May 2023, doi: 10.47852/bonviewJCCE3202720.
- [4] S. Garg, "Predictive Analytics and Auto Remediation using Artificial Intelligence and Machine learning in Cloud Computing Operations," *Int. J. Innov. Res. Eng. Multidiscip. Phys. Sci.*, vol. 7, no. 2, 2019, doi: 10.5281/zenodo.15362327.
- [5] Z. Wang, L. Xia, H. Yuan, R. S. Srinivasan, and X. Song, "Principles, research status, and prospects of feature engineering for data-driven building energy prediction: A comprehensive review," *J. Build. Eng.*, vol. 58, p. 105028, Oct. 2022, doi: 10.1016/j.jobe.2022.105028.
- [6] S. Thangavel, S. Srinivasan, S. B. V. Naga, and K. Narukulla, "Distributed Machine Learning for Big Data Analytics: Challenges, Architectures, and Optimizations," *Int. J. Artif. Intell. Data Sci. Mach. Learn.*, vol. 4, no. 3, pp. 18–30, Oct. 2023, doi: 10.63282/3050-9262.IJAIDSML-V4I3P103.
- [7] H.-N. Huang *et al.*, "Employing feature engineering strategies to improve the performance of machine learning algorithms on echocardiogram dataset," *Digit. Heal.*, vol. 9, Jan. 2023, doi: 10.1177/20552076231207589.
- [8] S. Garg, "AI/ML Driven Proactive Performance Monitoring, Resource Allocation and Effective Cost Management in SaaS Operations," *Int. J. Core Eng. Manag.*, vol. 6, no. 6, pp. 263–273, 2019.
- [9] R. Aurangzaib, W. Iqbal, M. Abdullah, F. Bukhari, F. Ullah, and A. Erradi, "Scalable Containerized Pipeline for Real-time Big Data Analytics," in *2022 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*, IEEE, Dec. 2022, pp. 25–32. doi: 10.1109/CloudCom55334.2022.00014.
- [10] A. Baluta, J. Mukherjee, and M. Litoiu, "Machine Learning based Interference Modelling in Cloud-Native Applications," in *Proceedings of the 2022 ACM/SPEC on International Conference on Performance Engineering*, 2022, pp. 125–132. doi: 10.1145/3489525.3511677.
- [11] P. Chandrashekar, "Advancements in Automated Incident Management : A Survey within Cloud-Native SRE (Site Reliability Engineering) Practices," *Int. J. Curr. Eng. Technol.*, vol. 13, no. 6, pp. 601–609, 2023, doi: 10.14741/ijcet/v.13.6.13.
- [12] C. Cheng, S. Li, and H. Ke, "Analysis on the Status of Big Data Processing Framework," in *2018 International Computers, Signals and Systems Conference (ICOMSSC)*, IEEE, Sep. 2018, pp. 794–799. doi: 10.1109/ICOMSSC45026.2018.8941875.
- [13] J. Liu, T. Zhu, Y. Zhang, and Z. Liu, "Parallel Particle Swarm Optimization Using Apache Beam," *Information*, vol. 13, no. 3, p. 119, Feb. 2022, doi: 10.3390/info13030119.
- [14] N. Mungoli, "Scalable, Distributed AI Frameworks: Leveraging Cloud Computing for Enhanced Deep Learning Performance and Efficiency," 2021, doi: 10.48550/arXiv.2304.13738.
- [15] A. Kushwaha, P. Pathak, and S. Gupta, "Review of optimize load balancing algorithms in cloud," *Int. J. Distrib. Cloud Comput.*, vol. 4, no. 2, pp. 1–9, 2016.
- [16] V. Rudel, P. Kienast, G. Vinogradov, P. Ganser, and T. Bergs, "Cloud-based process design in a digital twin framework with integrated and coupled technology models for blisk milling," *Front. Manuf. Technol.*, vol. 2, Dec. 2022, doi: 10.3389/fmtec.2022.1021029.

- [17] Y. Zhang and Z. Wang, "Feature Engineering and Model Optimization Based Classification Method for Network Intrusion Detection," *Appl. Sci.*, vol. 13, no. 16, p. 9363, Aug. 2023, doi: 10.3390/app13169363.
- [18] G. Sarraf and V. Pal, "Privacy-Preserving Data Processing in Cloud: From Homomorphic Encryption to Federated Analytics," *Int. J. Sci. Res. Comput. Sci. Eng. Inf. Technol.*, vol. 8, no. 2, pp. 735–706, 2022.
- [19] P. Chadha, "A Survey on Feature Engineering and Data Visualization for Result Predictions," *Int. J. Comput. Sci. Trends Technol.*, vol. 10, no. 4, pp. 104–111, 2022.
- [20] B. R. Cherukuri, "Microservices and containerization: Accelerating web development cycles," *World J. Adv. Res. Rev.*, vol. 6, no. 1, pp. 283–296, Apr. 2020, doi: 10.30574/wjarr.2020.6.1.0087.
- [21] P. Nama, "Integrating AI with cloud computing: A framework for scalable and intelligent data processing in distributed environments," *Int. J. Sci. Res. Arch.*, vol. 6, no. 2, pp. 280–291, Aug. 2022, doi: 10.30574/ijrsra.2022.6.2.0119.
- [22] K. Al-Barznji, "Big Data Processing Frameworks for Handling Huge Data Efficiencies and Challenges: A Survey," *Int. J. Data Sci. Big Data Anal.*, vol. 2, no. 1, pp. 1–9, May 2022, doi: 10.51483/IJDSBDA.2.1.2022.1-9.
- [23] R. Tandon and D. Patel, "Evolution of Microservices Patterns for Designing Hyper-Scalable Cloud-Native Architectures," *ESP J. Eng. Technol. Adv.*, vol. 1, no. 1, pp. 288–297, 2021, doi: 10.56472/25832646/JETA-VIIIIP131.
- [24] J. Kachhia, R. Natharani, and K. George, "Deep Learning Enhanced BCI Technology for 3D Printing," in *2020 11th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON)*, IEEE, Oct. 2020, pp. 0125–0130. doi: 10.1109/UEMCON51285.2020.9298124.
- [25] G. Toffetti, S. Brunner, M. Blöchlinger, J. Spillner, and T. Bohnert, "Self-managing cloud-native applications: Design, implementation, and experience," *Futur. Gener. Comput. Syst.*, vol. 72, 2016, doi: 10.1016/j.future.2016.09.002.
- [26] S. K. Chintagunta, "Survey of Containerization, Orchestration, and CI / CD Integration on DevOps in Modern Software Development," *Int. J. Curr. Eng. Technol.*, vol. 13, no. 6, pp. 610–618, 2023.
- [27] J. Lee *et al.*, "Comparative effectiveness of medical concept embedding for feature engineering in phenotyping," *JAMIA Open*, vol. 4, no. 2, pp. 1–12, Apr. 2021, doi: 10.1093/jamiaopen/ooab028.
- [28] T. Rawat and V. Khemchandani, "Feature Engineering (FE) Tools and Techniques for Better Classification Performance," *Int. J. Innov. Eng. Technol.*, vol. 8, no. 2, 2017, doi: 10.21172/ijiet.82.024.
- [29] C. Shan, Y. Xia, Y. Zhan, and J. Zhang, "KubeAdaptor: A docking framework for workflow containerization on Kubernetes," *Futur. Gener. Comput. Syst.*, vol. 148, pp. 584–599, Nov. 2023, doi: 10.1016/j.future.2023.06.022.
- [30] P. Olaya *et al.*, "Enabling Scalability in the Cloud for Scientific Workflows: An Earth Science Use Case," in *2023 IEEE 16th International Conference on Cloud Computing (CLOUD)*, IEEE, Jul. 2023, pp. 383–393. doi: 10.1109/CLOUD60044.2023.00052.
- [31] J. Shin, D. Arroyo, A. Tantawi, C. Wang, A. Youssef, and R. Nagi, "Cloud-native workflow scheduling using a hybrid priority rule and dynamic task parallelism," in *Proceedings of the 13th Symposium on Cloud Computing*, 2022, pp. 72–77. doi: 10.1145/3542929.3563495.
- [32] S. Henning and W. Hasselbring, "A configurable method for benchmarking scalability of cloud-native applications," *Empir. Softw. Eng.*, vol. 27, no. 6, p. 143, Nov. 2022, doi: 10.1007/s10664-022-10162-1.
- [33] C. Peña-Monferrer, R. Manson-Sawko, and V. Elisseev, "HPC-cloud native framework for concurrent simulation, analysis and visualization of CFD workflows," *Futur. Gener. Comput. Syst.*, vol. 123, pp. 14–23, Oct. 2021, doi: 10.1016/j.future.2021.04.008.
- [34] J. George *et al.*, "A Scalable and Cloud-Native Hyperparameter Tuning System," pp. 1–10, Jun. 2020.