

## Automated Feature Engineering Systems in Large-Scale Healthcare Data Environments

Velangani Divya Vardhan Kumar Bandi<sup>1</sup>

<sup>1</sup>Director AI/ML Engineering

Email: divyavardhanbandi@gmail.com,

Cite this paper as: Velangani Divya Vardhan Kumar Bandi (2024) Automated Feature Engineering Systems in Large-Scale Healthcare Data Environments. Journal of Neonatal Surgery, 13, 2127-2141

### ABSTRACT

A rapidly proliferating corpus of clinical research harnessing the power of machine learning has substantial implications for healthcare feature engineering. As a broad umbrella encompassing data preprocessing, quality control, transformation, and generation, feature engineering addresses a major bottleneck in the production of predictive models. Automated feature engineering systems are increasingly deployed at scale to meet the challenges of generating the vast quantity of predictive features necessary for successful, generalizable, and clinically useful machine learning systems. Such clinical feature engineering systems produce features that are applied in a predictive setting after the fact and not explicitly linked to clinical care, but nevertheless involve substantial risk. A principled examination of a clinical feature engineering system can be framed in terms of six core components: data governance; compliance with privacy and regulatory constraints; appropriate validation and prospective evaluation; consideration of data biases; the use of effective data-quality and preprocessing pipelines; and sound candidate feature generation, scoring, and selection strategies. Health systems typically possess an assemblage of rich and diverse, yet underutilized, information with the potential to contribute meaningfully to clinical prediction problems. Electronic health record (EHR) data, comprising clinical notes, laboratory values, medication orders, and procedure codes; over a decade's worth of length and width dataset and point-of-care laboratory test results; continuous biosensor measurements; DNA sequencing data; biomarkers derived from imaging; and drug compounds targeting genotypes provide raw material for hundreds of prediction problems in diverse specialties. However, machine learning in healthcare exhibits a stunning lack of reproducibility: many predictive models fail to retain their accuracy in different cohorts, and those that do are seldom incorporated into routine clinical care. A significant bottleneck underlying this failure lies with the feature engineering step...

**Keywords:** *Clinical Feature Engineering, Automated Feature Engineering Systems, Healthcare Machine Learning, Predictive Model Pipelines, Electronic Health Records (EHR), Multimodal Clinical Data, Data Governance in Healthcare AI, Privacy and Regulatory Compliance, Prospective Model Validation, Bias-Aware Feature Design, Data Quality and Preprocessing Pipelines, Candidate Feature Generation and Selection, Feature Scoring Strategies, Reproducible Clinical ML, Biosensor and Time-Series Data, Imaging-Derived Biomarkers, Genomic and Sequencing Features, Point-of-Care Laboratory Data, Clinical Model Generalization, Production-Grade Healthcare AI.*

### 1. INTRODUCTION

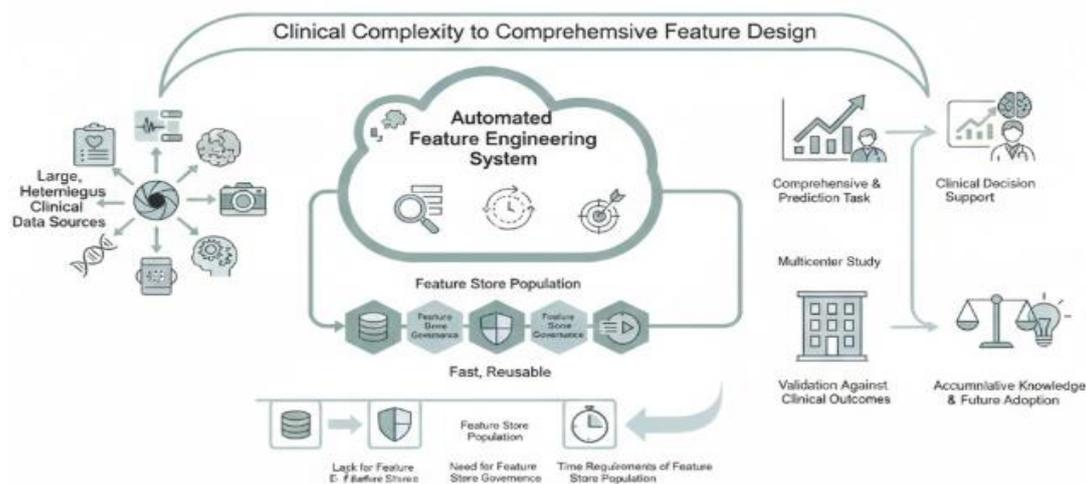
Feature engineering is a challenging, yet critical, part of the supervised machine learning pipeline. While there are many exciting, new ML models that have demonstrated their success in other fields, such as language, vision, or games, they rely on large, curated datasets that require a careful selection, crafting, and engineering of features. Following a similar approach in healthcare is very difficult, due mainly to the many different sources of data that are used, the fact that they originate from different hospitals and that they must be updated each time a new prospective study is proposed. Because of these challenges, such studies can take months or even years before a model has been built, tested, and validated. Ideally, careful consideration of the features and data sources should counterbalance that, improving predictive performance, enhancing generalization, or leading to the discovery of new risk factors. At least from an algorithmic perspective, it would be sufficient to automate feature engineering and treat it as a black box. Automating the generation of predictive features in healthcare then pushes these methods one step further.

Healthcare feature engineering is the process of gathering, selecting, and preparing features from one or more healthcare data sources, in particular Electronic Health Records (EHR) databases. Such data sources contain a wealth of information on patients, including their medical history, medications, and diagnostics, which can be used to predict future clinical outcomes. However, these sources vary from site to site and evolve over time, introducing heterogeneity and uncertainty into the dataset. Therefore, feature generation consists mainly of three stages: generating a source of feature candidates by applying domain

knowledge and heuristics, deploying data preparation and cleaning pipelines to ensure that the features produced can then be used safely in a future model, and scoring the candidates with an automated or manual process to finally select the best-ranked features. Past success of the first stage in using domain knowledge and heuristics to mine predictive subgraphs with temporal constraints has demonstrated the potential of feature generation methods in healthcare

### 1.1. Overview of Healthcare Feature Engineering

Healthcare feature engineering encompasses the design of features from large, heterogeneous, longitudinal clinical data sources—particularly electronic health records, but also imaging, genomic, sensor, and textual data—to support clinical decision-making.



**Fig 1: Scaling Clinical Intelligence: Automated Feature Engineering and Governed Feature Stores for Multimodal, Reproducible Healthcare AI**

It aims to produce a comprehensive set of clinically relevant, time-varying features for a targeted patient population that maximally inform the prediction task, integrating all data sources available for those patients. Clinical complexity makes comprehensive feature design highly challenging, motivating automated feature engineering systems that mirror the pipeline of data scientists with rich domain knowledge.

Two features or sets of features are rarely used across studies, limiting the accumulative knowledge-building typical in research. While proposed feature sets are often extensive, they typically focus on a single domain (e.g., EHR diseases, meds, lab results) without standardized integration across domains. The reasons are threefold: lack of a feature store, the need for feature store governance, and the time requirements of feature store population. An automated feature engineering system fulfills these needs and enables fast, comprehensive, and reusable feature engineering for a target study population. Applying such a system in a multicenter study provides insights on usage patterns and enables validation against clinical outcomes, supporting future adoption of broader, prospective feature design pipelines.

## 2. Foundations of Feature Engineering in Healthcare

Considerable amounts of semi-structured, unstructured, and high-dimensional medical data have been amassed over the past few decades. Prominent data sources include electronic health records (EHRs), medical imaging and diagnostic reports, genomic sequencing data, biosensors, mobile health applications, and wearable sensors. These modalities can complement clinical data, provide additional insights about a patient's condition, or even replace clinical assessments such as the ejection fraction score for heart disease. Despite their apparent utility in clinical decision-making, these auxiliary datasets are still seldomly used in supervised machine learning tasks.

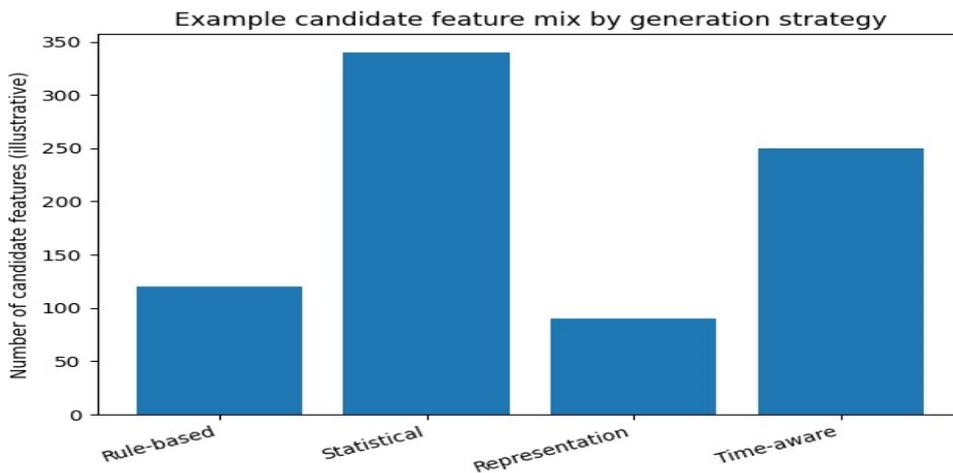
Within the context of supervised machine learning, a feature is defined as a measurable characteristic of a phenomenon being observed. Feature stores serve as centralized repositories for storing and serving data features for use in machine learning models. These features are engineered and made available for use in supervised learning tasks, with properties that allow for easy discovery, management, and use across projects. Features store metadata also provide a source of potential value by tracking data provenance. Well-maintained feature stores can significantly reduce the time and effort to develop ML models

and increase reproducibility by allowing for easy retrieval of the same feature set.

### 2.1. Clinical Relevance and Data Modalities

Healthcare machine learning relies on diverse clinical datasets and requires features with medical interpretations that are relevant to clinical decision-making. The most common data modalities in healthcare are electronic health records, medical images, genomic data, and data from wearable and stationary sensors. The data sources are inherently different in their content and the dimensions they measure which is further complicated by the vast amount of heterogeneous information that is aggregated across time. In addition, most of the data sources available for ML are generated during the normal course of clinical practice and when other scientific investigations take place which raises challenges for data quality. Proper mappings must enable the application of ML techniques on the collected data sources.

In healthcare environments, data sources are most frequently accessed in keep-it-simple strategies, whether through the data lake/exploitation or the hospital's multiple data products. Authors focus on the first strategy, where a clinic uses the data lake and the other clinics make use of the data distributed in a subsequent event. When taking this approach, data quality is still an issue and authors combined information cleaning and data harmonization cycle, applied on retrospective healthcare ML applications, with validation in cross-site frameworks. The feature engineering process must then take the complete process into account but be divided in two distinct machines: one for data quality assurance and the second for feature generation.



**Fig 2: Data Quality and Preprocessing Workflow: Cleaning, Normalization, Imputation, De-Duplication, and Harmonization in Clinical Feature Pipelines**

#### Equation 1) Data Quality & Preprocessing Equations (Clean → Normalize → Impute → Deduplicate → Harmonize)

##### 1.1 Cleaning As Constraint Projection (Validity Rules)

Let raw measurement be  $x$ . Suppose “plausible clinical bounds” for a lab are  $[L, U]$ . A standard cleaning operation is clipping (winsorization):

If  $x < L$ , set  $x' = L$

If  $L \leq x \leq U$ , set  $x' = x$

If  $x > U$ , set  $x' = U$

Compactly:

$$x' = \min(\max(x, L), U)$$

##### 1.2 Normalization

**(A) Min-max normalization** (map to  $[0,1]$ )

Given a feature column  $x_1, \dots, x_n$ :

Compute  $x_{\min} = \min_i x_i, x_{\max} = \max_i x_i$

For each value:

$$Z_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

**(B) Z-score standardization** (mean 0, variance 1)

Mean:

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

Population std:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

Standardize:

$$z_i = \frac{x_i - \mu}{\sigma}$$

### 1.3 Imputation (missing values)

Let  $x_i$  be missing for some indices.

#### Median imputation

Compute median  $m = \text{median}(\{x_i: x_i \text{ observed}\})$

Replace missing values:

$$x_i^{\text{imp}} = \begin{cases} x_i, & x_i \text{ observed} \\ m, & x_i \text{ missing} \end{cases}$$

**Mean imputation** is the same form with  $m \leftarrow \mu$ .

**Linear interpolation** (for time series at times  $t_k$ )  
 If  $x(t_a)$  and  $x(t_b)$  are observed and  $t_a < t < t_b$  is missing:

$$x(t) = x(t_a) + \frac{t - t_a}{t_b - t_a} (x(t_b) - x(t_a))$$

## 2.2. Defining Features, Features Stores, and Lineage

A feature is a single measurable property or characteristic of a phenomenon. In a healthcare context, clinical features refer to attributes or constructs captured through diagnostic or prognostic procedures, exam findings, or tests. Laboratory test results represent features measured at a point in time. Genomic expression analysis examines hundreds of thousands of features simultaneously, each referring to expression levels of specific genes. Medical imaging is fundamentally about measuring features, and different imaging modalities extract different features relevant for specific disease processes. A feature store is a repository for features, usually organized by their predictive task. Feature stores facilitate discovery and reuse. Automated Feature Engineering in a large healthcare environment implies resilient and scalable software products. Features are gathered, modeled, stored, and tested for algorithmic fitness in a systemic way.

Maintaining feature lineage—tracking how features were created and the precise data (origin, version, state) and processing pipelines used to generate them—greatly aids collaboration, improves reproducibility, addresses quality issues, and enhances compliance with regulatory requirements. Lineage tracking typically follows a hierarchical approach: features combine datasets, and datasets depend on raw data and other derived datasets. The entire history of the computation underlying a clinical feature can be captured and stored with metadata about the feature. Individual nodes in the lineage can be attributed and versioned, can record the health of the lineage, and can be used to trigger de-duplication, re-creation, and recalibration. Individual features may have governance models.

**Table 1: Sample Deduplicated Clinical Feature Records for Patients**

patient_id	time_hr	creatinine_after_dedup	age
101	0	1.1	55
102	0		70
103	0	0.6	40
104	0	9.8	62
105	0	1.4	33

**Automated Feature Engineering: Techniques and Architectures**

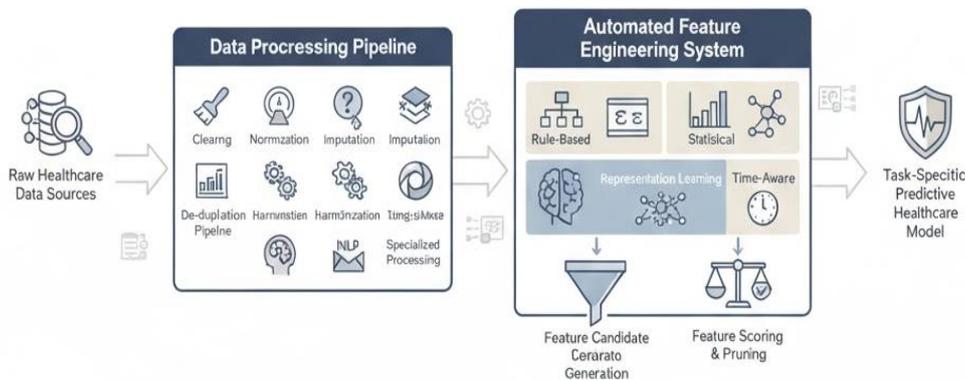
Automated feature engineering encompasses four broad categories: rule-based generation, statistical feature generation, representation learning, and the generation of time-aware features. Each category may apply to a single feature or groups of features. Rule-based features may be generated with any suitable method (e.g., user or crowd-sourced rules, templates, philosophy-based ontologies). Statistical feature generation constitutes the application of functions to selected groups of features across any user-defined window (a.k.a., moving or sliding window). Representation learning encapsulates the suite of (semi-)unsupervised techniques for learning representations from large data sets. Finally, time-aware features are those whose values depend explicitly on time, areas of time, or other temporal characteristics.

With a library of candidate regions of interest generated by detection, a data provenance strategy designed to facilitate tracking and scoring candidate region features, and a functional pipeline in place to clean and prepare the region of interest features, the remaining open questions may be framed along two dimensions: the data quality aspects of preprocessing and the implementation of a pipeline capable of effecting (1) common operations such as input data cleaning, normalization, de-duplication, and temporal alignment across sites, and (2) specific target definitions calling for advanced harmonization (i.e., creation of features incorporating data from more than one site).

**3.1. Feature Generation Strategies**

A variety of techniques may be applied to automate feature generation from healthcare data. Literature identifies methods that can be categorized as rule-based approaches, statistical methods of feature generation, data-centric representation learning methods that facilitate the construction of task-specific features directly from the data, and time-aware feature generation techniques. Regardless of the specific techniques utilized, candidates must be scored—distinguishing between features worth pursuit and those that should be automatically pruned—and scored candidates subsequently selected.

Conceptually, healthcare features are derived from the processing of raw data: either a transformation is applied to the underlying data or a set of other features is used to train a predictive model for a new feature. Data preprocessing methods, including cleaning, normalization, imputation, de-duplication, harmonization across different data-collection sites, and domain-specific preprocessing workflows (e.g., image and natural-language processing), can be implemented in a standalone data-quality pipeline prior to candidate feature engineering.



**Fig 3: Automated Feature Engineering Systems in Healthcare: A Multimodal Pipeline for Candidate Generation, Scoring, and Clinical Validation**

### 3.2. Data Quality and Preprocessing Pipelines

A second pipeline addresses data quality and feature preprocessing tasks. Cleaning removes invalid, improbable, and contradictory values. Normalization scales values to a common range while preserving semantic meaning. Imputation fills in missing values through median/mode, value interpolation, predictive models, or generative methods. De-duplication detects and reconciles duplicate records for the same entity at the same point in time. Harmonization reconciles semantically equivalent but contextually different records across data sources. Feature-specific preprocessing handles transformations, discretization, binning, embeddings, dimensionality reduction, and algorithm-specific transformations such as one-hot encoding for categorical variables.

Maintaining a Feature Store to centralize and standardize features while supporting lineage tracking increases reproducibility and transparency. A feature is an individual measurable characteristic or property of a phenomenon being observed. Automated feature engineering seeks to reduce the manual effort that data scientists spend on feature creation for machine learning models by developing tools, algorithms, and system support. Three automated steps are candidate generation, scoring, and selection.

### 4. Systems for Large-Scale Healthcare Environments

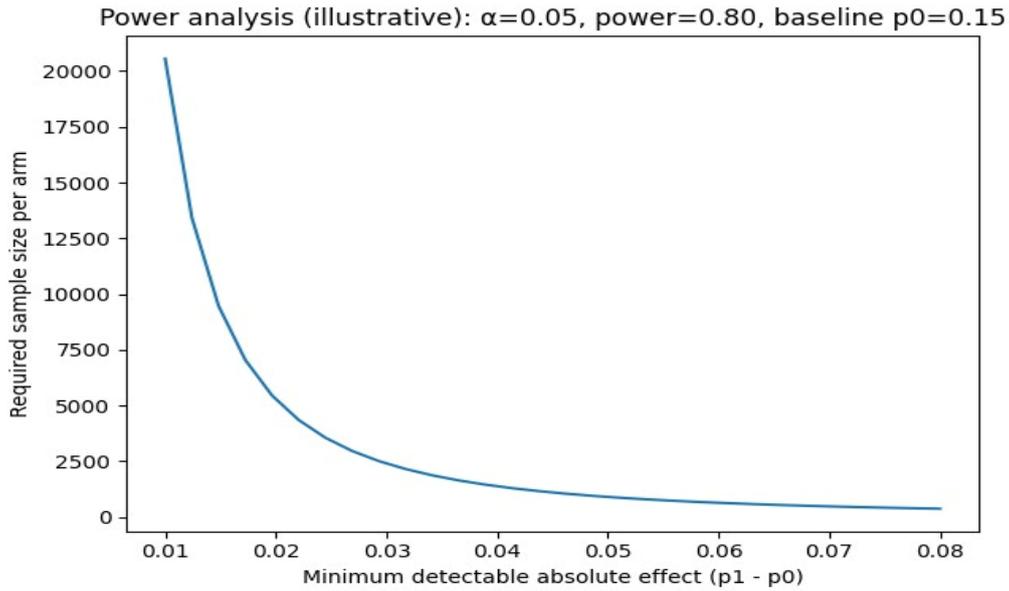
Data governance encompasses access control and logging procedures, de-identification methods (pseudonymization, k-anonymization, data masking), and assessment of trusted data-sharing frameworks (secure multi-party computation, trusted third-party-computation), allowing multiple parties to analyze sensitive data without revealing source information. Because clinical data often span multiple sites for many patients, harmonization needs to go beyond cleaning and normalizing the data to feature generation and candidate scoring processes. Data privacy and security must be firmly established to assure compliance with healthcare regulations. With reference to the United States, Australia, and Europe, regulatory oversight includes not only the Health Insurance Portability and Accountability Act and the General Data Protection Regulation but explicit links to Food and Drug Administration regulation of medical devices and drug trials, and considerations for limited patient consent and health-information-sharing agreements.

Evaluation and validation in retrospective studies focus on generalization across sites and populations, while prospective assessment establishes the safety and clinical effectiveness of a deployed system. Beyond receiver-operating-characteristic-area-under-the-curve scores, true clinical utility hinges on statistical power and, hence, the extent and timing of study recruitment. Post-validation management tasks—monitoring for adverse behavior, testing statistical hypotheses of changed performance, and executing planned cross-validation—respond to the open-endedness of clinical studies through exploratory analysis and pondering the root causes of any emerge performance shifts.

#### 4.1. Data Governance, Privacy, and Security

A comprehensive data governance framework assures that data remains secure, private, and confidential even when used across many institutions in a distributed, collaborative environment. Mechanisms for governing access control to data must be granular and fine-tuned to give the right user the appropriate levels of access for their current task. Audit logging keeps track of data usage; it allows tracking of the frequency and patterns of access, thereby revealing suspicious and potentially malicious activity. When raw clinical datasets contain sensitive personally identifiable information (PII), data must undergo de-identification before sharing. Sharing across sites using secure multi-party computation (MPC) allows complex joint analysis without disclosing PII, but it tends to be slow and users need specialized programming skills.

Healthcare data is sensitive by nature. Leaving unmonitored access to sensitive data across multiple external institutions is risky for the data owners, creating privacy and confidentiality issues. Compliance with healthcare privacy regulations, such as the Health Insurance Portability and Accountability Act (HIPAA) in the United States or the General Data Protection Regulation (GDPR) in Europe, is essential in a healthcare data governance plan. Data governance metrics involve identifying what data has been processed and how. Data retention policies describe how long each dataset should be kept and how it will be securely disposed of afterwards.



**Fig 4: Time-Aware Statistical Feature Engineering: Moving Window Aggregations, Variance Estimation, and Trend Extraction Over Longitudinal Patient Data**

**Equation 2) Time-Aware & Statistical Feature Generation (Moving Windows)**

Let patient  $p$  have measurements  $\{(t_j, x_j)\}$ . For a window length  $W$ , define the set in the window ending at time  $t$ :

$$\mathcal{J}(t) = \{j: t - W < t_j \leq t\}$$

**2.1 Window mean**

$$\text{mean}_W(t) = \frac{1}{|\mathcal{J}(t)|} \sum_{j \in \mathcal{J}(t)} x_j$$

**2.2 Window variance**

Window mean  $\bar{x}_W(t)$  as above

Then

$$\text{var}_W(t) = \frac{1}{|\mathcal{J}(t)|} \sum_{j \in \mathcal{J}(t)} (x_j - \bar{x}_W(t))^2$$

**2.3 Trend (slope) in a window (linear regression)**

Fit  $x \approx \beta_0 + \beta_1 t$  for points in  $\mathcal{J}(t)$ .

Let window points be  $(t_j, x_j)$  for  $j \in \mathcal{J}(t)$ . Define:

$$\bar{t} = \frac{1}{m} \sum t_j, \quad \bar{x} = \frac{1}{m} \sum x_j, \quad m = |\mathcal{J}(t)|$$

Then the least squares slope is:

$$\beta_1 = \frac{\sum (t_j - \bar{t})(x_j - \bar{x})}{\sum (t_j - \bar{t})^2}$$

This produces a time-aware “trend” feature (e.g., creatinine rising over 48h).

**4.2. Compliance With Healthcare Regulations**

Frameworks mapping automated feature engineering into three classes of health care regulations include the HIPAA Privacy and Security Rules, the European Union General Data Protection Regulation (GDPR), and regulatory considerations from the U.S. Food and Drug Administration (FDA). Retention time is another important regulatory factor determining the

frequency data are updated and the prospective use of a feature, such as when it is included in a data retention policy for automated and continuous feature operation. Consent, data use agreements, and contractual requirements from institutional review boards (IRBs) and endpoint-contributing data owners are considered for Trans-Union and Trans-Atlantic multimodality healthcare feature repositories.

**Table 2: Automated Feature Ranking Using Predictive Strength, Stability, and Leakage Risk**

feature	predictive_score	stability_across_sites	leakage_risk
mean_creatinine_24h	0.42	0.9	0.1
trend_creatinine_48h	0.38	0.76	0.15
med_count_7d	0.31	0.88	0.08
note_embedding_dim17	0.35	0.7	0.2
icu_bed_occupancy_site	0.12	0.55	0.05

## 5. Evaluation, Validation, and Deployment

Adequate evaluation is critical to determine if automated features are safe for clinical use and capable of performing as well as or better than hand-engineered alternatives. Retrospective validation of candidate features in electronic healthcare record data through offline modelling provides an initial assessment of practical utility. This validation can encompass standard validation techniques used in machine learning, including train-test splits, cross-site validation for assessing generalizability, performance calibration, and error analysis. Candidates are then deployed into prospective clinical settings in an automated manner with pipelines that monitor for scoring errors and track conciliation. A/B testing of predictive and prognostic models (and potentially risk stratification scores) in live clinical settings — as facilitated by the availability of a feature store — provides an opportunity to assess value of features, ideally against clinical endpoints such as in-hospital mortality, ICU admissions, or length of stay.

Ensure cross-site study population sizes are suitably large and that analyses can be stratified when necessary to cover potential differences in feature efficacy. Allocating sufficient time until results become available is critical. Ideally, deploy interventions with sufficiently low expected effect size, otherwise statistical power rapidly becomes overly demanding. When modelling the likely size of the intervention effect, consider the time-horizon and relevant data epochs. Upon deployment, design the statistical tests with the necessary  $\alpha$  and  $\beta$  values, determine the minimum difference of interest, and calculate the necessary sample size. Ensure suitable monitoring mechanisms are in place to detect any anomalies after intervention. Finally, ascertain patient safety through trigger systems capable of identifying deleterious outcomes.

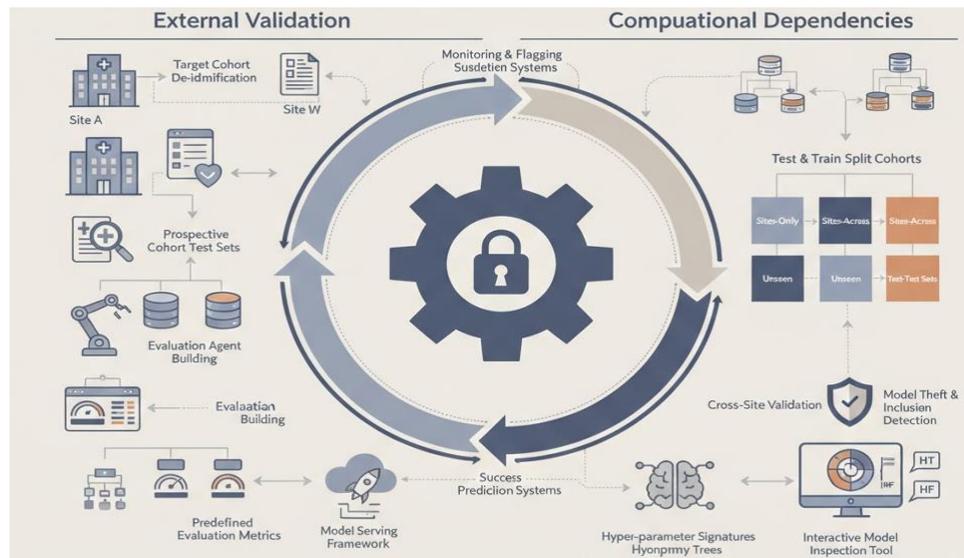
### 5.1. Offline Validation and Cross-Site Generalization

Worker performance, in terms of the AI model serving as a sub-component, must therefore be assessed in some detail. Macro-level validation targets for each pair of sites—e.g. a HF test set specifically for HF patients in sites A and B—must be established and monitored at deployment time. Critical metrics are calibrated, and candidate agents flagged according to their error analysis.

A standard external validation pipeline has also been established, including the generation of separate test sets from prospective cohorts (i.e. operating on data unseen by the learning agents). Initial preparations of the target cohorts, including de-identification, are performed by Site A. Evaluation agents are then built for the cohorts, with performance estimates returned to Site A. A predefined set of evaluation metrics is applied that enable dimension-based checks and considerations of medical validity for the evaluation results. A second pipeline applies the same checks for HT models, with agents routed through a model serving framework to facilitate the process.

At the other end, the computational dependencies of supervised learning models can also be scrutinized. By tracking feature use, test and train split colors are applied to pre-train/test cohorts, enabling a clear view of sites-only, sites-across, unseen, and test-test test sets. Cross-site validation by test split color further allows rapid scans of known strength for model theft and inclusion detection—enough to inform the use of success prediction systems.

Dominant logical mechanisms and hyper-parameter signatures within trained models can finally be parsed, enabling the recovery of hyponymy trees for an interactive model inspection tool. The visualization ensures that summarization can extend to modelling efforts across HT and HF elements.



**Fig 5: Cross-Site Clinical Validation: A Framework for Multi-Cohort Evaluation, Computational Dependency Tracking, and Interactive Model Inspection**

### 5.2. Prospective Evaluation and A/B Testing In Clinical Settings

After offline validation, a deployment pipeline is used to move model-dependent features into a separate features store, followed by prospective evaluation in a clinical setting. Live traffic is assigned to candidate and baseline features following randomized assignment, with discernible clinical endpoints and an a priori power analysis determining sample size. Monitoring systems track commencement of the trial, and safety checks are established to stop the trial if the trial-determining feature fails to have a considerable effect in the direction signaled by the trials-determining model.

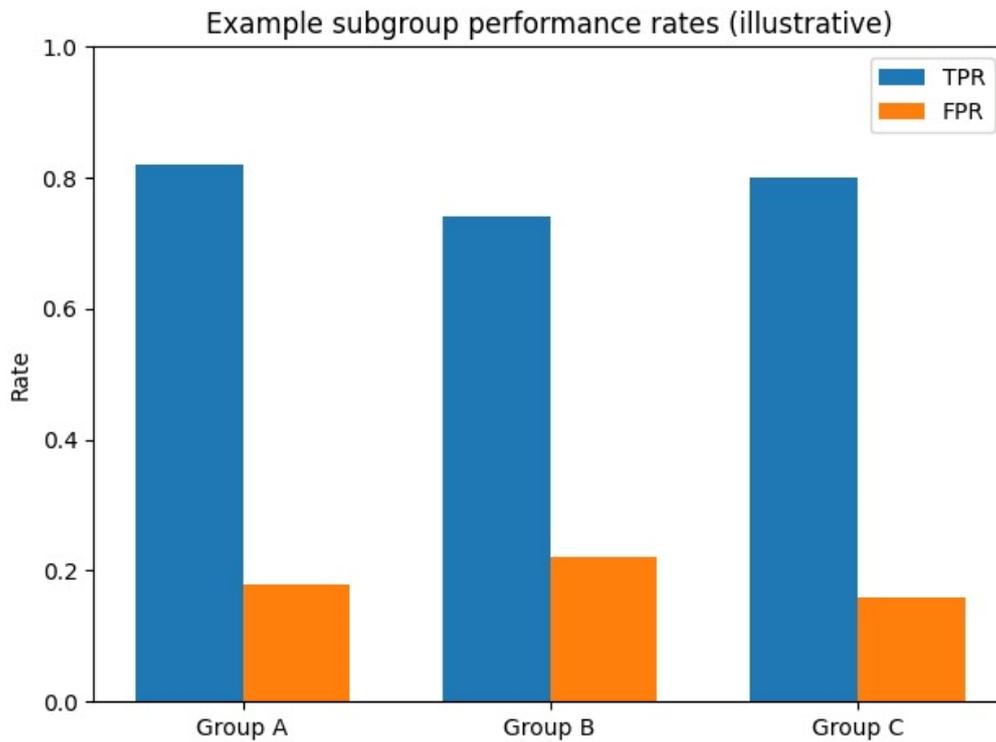
Prospective evaluation using A/B testing methodology is a natural next step for the validation of candidate features once offline requirements are met. This procedure permits the quantification of the clinical effect of the newly supplied model-dependent feature, retaining the benefits of assignment to features in production as feature-dependent active learning. The flexibility afforded by reduced data-collection costs, high accuracy, and substantial effort in de-biasing of the model-determining features permits formal A/B testing.

## 6. Challenges and Risks in Automated Feature Engineering

**Automated Feature Engineering Systems in Large-Scale Healthcare Data Environments:** High-dimensional outcome-sensitive feature engineering critically influences the outcome of machine learning, yet few methods exist. In clinical settings, high-dimensional EHR and imaging data may contain features that are costly or risky to evaluate prospectively, necessitating automated feature engineering. Data governance, privacy, and security must satisfy industry standards; automated feature generation and risk mitigation tools should be tailored to healthcare’s unique characteristics.

**Data Bias, Fairness, and Equity Implications:** Automated feature engineering can inadvertently produce problematic models that encode racial, gender, or economic bias. Systemically, analysis can estimate fairness-related model properties across protected attributes, examine performance on subgroups, and compute group-wise fairness metrics; remediation strategies, if needed, include sample stratification and in-processing adjustment. From a risk perspective, approaches in healthcare data have shown that building a clinical pathway for an outcome of interest allows testing for bias as well as links to social determinants of health, making bias, fairness, and equity implications visible.

**Privacy-Preserving Feature Generation:** Scalable healthcare feature engineering leverages data of unknown ownership via core components for privacy-sensitive training and feature generation. Federated learning, a distributed machine learning approach, trains models closest to the raw data, sidestepping raw data transfer and ownership concerns. Privacy accounting and differential privacy ensure model release and evaluation satisfy privacy requirements. Within an automated pipeline, potential feature generators are ranked by privacy risk, and resource-aware differential privacy correctly weights real and synthetic features. Reported privacy risk estimates promote stakeholder trust; clearly articulated costs lower adversarial investment.



**Fig 6: Candidate Feature Scoring and Selection: Predictive Strength, Stability Across Sites, Leakage Risk Penalization, and Combined Ranking Strategy**

**Equation 3) Candidate feature scoring & selection (a basic, complete math formulation)**

Let there be candidate features  $f_1, \dots, f_M$ . Define three common scoring ingredients:

**3.1 Predictive strength (example: correlation with outcome)**

For numeric outcome  $y$ , Pearson correlation:

$$\text{Means: } \bar{f} = \frac{1}{n} \sum f_i, \bar{y} = \frac{1}{n} \sum y_i$$

$$\text{Covariance: } \text{cov}(f, y) = \frac{1}{n} \sum (f_i - \bar{f})(y_i - \bar{y})$$

Std devs:  $\sigma_f, \sigma_y$

Correlation:

$$r = \frac{\text{cov}(f, y)}{\sigma_f \sigma_y}$$

**3.2 Cross-site stability (example: variance of effect across sites)**

If you estimate a site-specific effect  $\theta_s$  (e.g., AUC lift or coefficient), stability can be inverse dispersion:

$$\text{stability} = \frac{1}{1 + \text{Var}_s(\theta_s)}$$

**3.3 Leakage / risk penalty**

If leakage score  $L(f) \in [0,1]$  (higher worse), then penalize.

**3.4 Combined ranking score**

A simple complete selection score:

$$S(f) = w_1 \cdot \text{predictive}(f) + w_2 \cdot \text{stability}(f) - w_3 \cdot \text{leakage}(f)$$

Select top- $K$ :

$$\{f\}_{\text{selected}} = \arg \text{top}KS(f)_{f \in \{f_1, \dots, f_M\}}$$

### 6.1. Data Bias, Fairness, and Equity Implications

Multiple sources of bias may lead to datasets that fail to adequately represent entire segments of the healthcare population, creating issues of predictive fairness and inequity. From an equity perspective, it is critical to identify the potential sources and forms of bias in automated feature engineering systems, as well as their implications for population-specific predictions. Systematic methods, such as those by Barocas et al., can be employed to investigate and characterize bias in both the data and resulting deployed model. The key steps are as follows:

1. Consider important attribution questions: Who or what is the object of concern? (e.g., individuals, groups, models, predictions) Who or what is responsible for the outcome? (e.g., individuals, groups represented in the data, models) Who or what is affected by the outcome? (e.g., individuals, groups, society) Is the outcome rendered unjust? If so, how? (e.g., inequitable, unfair, unworthy)
2. Identify potential sources of bias in data: Bias can stem from issues relating to data generation, collection, integration, modeling/engineering, evaluation/selection, and actual deployment.
3. Diagnose and measure bias: Precision diagnostics reveal underrepresented subpopulations and their associated metrics, and fairness metrics expose whether the model merits concern from a fairness perspective.
4. Mitigate detected imbalance: Attention needs to be devoted to addressing detected concerns, possibly employing stratification by risk or tailoring the model in other ways.

### 6.2. Privacy-Preserving Feature Generation

Privacy is paramount in healthcare, where sensitive patient information is collected, analyzed, and shared. The creation of features that rely on personal data must therefore be governed by data privacy and statistical disclosure risk principles to minimize the risk of patient identification. For cross-institutional sharing, privacy constraints often preclude direct data sharing and therefore necessitate special computational techniques that allow model training without direct access to underlying data. Federated learning, for example, enables decentralized ML wherein the algorithm is trained across multiple parties without exchanging the data used to build locally updated models. The final model is obtained by aggregating the locally trained models while meeting regulatory data use agreements. Other techniques, such as building a differentially private version of the features that aggregate information across a subset of users or designating a trusted data steward to manage the model training process, help ensure that direct or indirect identification is highly improbable. Therefore, a careful assessment of the risk of privacy violation and the definition of rigorous privacy budgets, together with differential privacy techniques, can facilitate privacy-compliant feature generation.

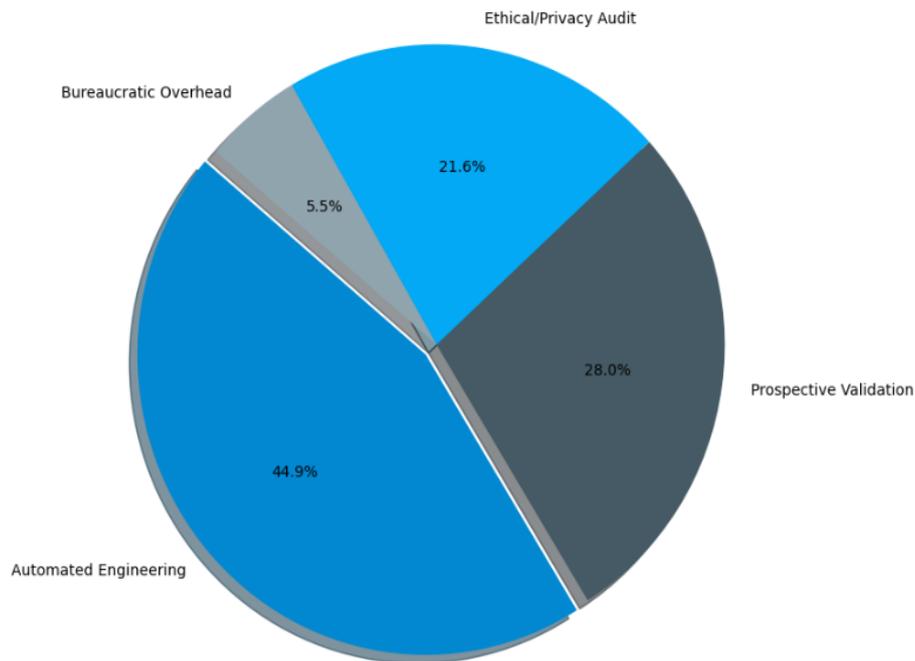
The importance of data privacy considerations increases when developing models capable of automated feature generation, such as GNNs or VAEs. Privacy-preserving training of these models has garnered increasing interest, driven by the success of ML and the associated imperative of having access to large amounts of data. Previous work has proposed a federated learning methodology for training GNN models in a healthcare setting, with focus on monitoring heterogeneous disease progressions via subgraph-based learning. Additionally, several methods for privacy-preserving generative modeling using VAE architectures have recently been proposed. Thus, proper design considerations can enable automated feature generation for highly sensitive datasets.

## 7. CONCLUSION

Automated feature engineering for predictive models trained on healthcare data is crucial because of the associated economic and human costs of healthcare data science projects. Feature stores specifically designed for large clinical datasets enable the sharing, quality control, and systematic generation of data quality-enhancing features, rules, and cleaning procedures among data scientists and stakeholders. A mandated feature store addresses the challenge of ensuring that all services utilize optimal rules and features, but at the expense of introducing a bureaucratic overhead.

Herein, the central role of features in predictive models for clinical data pipelines, the techniques and architectures for automation, and the specialized requirements for large healthcare data environments are substantiated. Potential sources of bias during automated feature generation and the implications of such biases are examined. Frameworks for prospective validations in clinical settings that go beyond model performance and for the ethical generation of features given privacy concerns are explored. Finally, avenues for resolving open methodological questions and translating automated feature engineering systems into production are identified.

## Feature Science Cost Management



**Fig 7: Feature Science Cost Management**

### 7.1. Summary and Future Directions in Healthcare Feature Engineering

Automated Feature Engineering is a crucial yet largely unexplored field in healthcare. Key technical aspects of Automated Feature Engineering Systems in Large-Scale Healthcare Environments are identified and characterized, encompassing feature generation strategies; procedures for assessing data quality and implementing preprocessing pipelines; requirements for data governance, privacy, and security; how the systems address HIPAA, GDPR, and FDA regulations; mechanisms for offline validation; and a framework for prospective evaluation and A/B testing in clinical settings.

Although Automated Feature Engineering holds considerable promise across multiple application domains, including healthcare, further exploration is necessary. For example, little is known about data bias, fairness, and equity implications; a novel set of heuristic fairness metrics has recently been proposed, yet practical approaches for bias mitigation remain scarce. Applying novel privacy-preserving techniques such as secure multiparty computation, federated learning, and differential privacy is another promising avenue, especially given heightened attention to data privacy and protection. Accelerating the application of Automated Feature Engineering in healthcare would ultimately enhance the safety, efficacy, and equity of predictive models, thereby addressing the pressing need for scalable solutions and enabling artificial intelligence to realize its full potential across healthcare and other domains.

### REFERENCES

- [1] Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D. G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., & Zheng, X. (2016). TensorFlow: A system for large-scale machine learning. In Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI) (pp. 265–283). USENIX.
- [2] Gottimukkala, V. R. R. (2023). Privacy-Preserving Machine Learning Models for Transaction Monitoring in Global Banking Networks. *International Journal of Finance (IJFIN)-ABDC Journal Quality List*, 36(6), 633-652.
- [3] Akhtar, A., Khan, M., & Nazir, S. (2021). Industrial anomaly detection: A survey of methods and applications. *Computers & Industrial Engineering*, 158, 107377.
- [4] IT Integration and Cloud-Based Analytics for Managing Unclaimed Property and Public Revenue. (2024). *MSW Management Journal*, 34(2), 1228-1248.
- [5] Alur, R., & Dill, D. L. (1994). A theory of timed automata. *Theoretical Computer Science*, 126(2), 183–235.
- [6] Angelopoulos, C. M., Nikolettseas, S., & Patroumpa, D. (2020). Edge computing in the Industrial Internet of Things: A survey. *IEEE Internet of Things Journal*, 7(10), 10665–10682.
- [7] Agentic AI in Data Pipelines: Self Optimizing Systems for Continuous Data Quality, Performance and Governance. (2024). *American Data Science Journal for Advanced Computations (ADSJAC) ISSN: 3067-4166*, 2(1).

- [8] Babiceanu, R. F., & Seker, R. (2016). Big Data and virtualization for manufacturing cyber-physical systems: A survey of the current status and future outlook. *Computers in Industry*, 81, 128–137.
- [9] Meda, R. (2024). Agentic AI in Multi-Tiered Paint Supply Chains: A Case Study on Efficiency and Responsiveness. *Journal of Computational Analysis and Applications (JoCAAA)*, 33(08), 3994-4015.
- [10] Bagheri, B., Yang, S., Kao, H.-A., & Lee, J. (2015). Cyber-physical systems architecture for self-aware machines in Industry 4.0 environment. *IFAC-PapersOnLine*, 48(3), 1622–1627.
- [11] Nagabhyru, K. C. (2024). Data Engineering in the Age of Large Language Models: Transforming Data Access, Curation, and Enterprise Interpretation. *Computer Fraud and Security*.
- [12] Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (pp. 610–623). ACM.
- [13] Kolla, S. H. (2024). RETRIEVAL-AUGMENTED GENERATION WITH SMALL LLMs FOR KNOWLEDGE-DRIVEN DECISION AUTOMATION IN ENTERPRISE SERVICE PLATFORMS. *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*, 15(3), 476–486.
- [14] Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13, 281–305.
- [15] Aitha, A. R. (2024). Generative AI-Powered Fraud Detection in Workers' Compensation: A DevOps-Based Multi-Cloud Architecture Leveraging, Deep Learning, and Explainable AI. *Deep Learning, and Explainable AI* (July 26, 2024).
- [16] Bonawitz, K., Ivanov, V., Kreuter, B., Marcedone, A., McMahan, H. B., Patel, S., Ramage, D., Segal, A., & Seth, K. (2017). Practical secure aggregation for privacy-preserving machine learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security* (pp. 1175–1191). ACM.
- [17] Bosch, J. (2018). Speed, data, and ecosystems: The future of software engineering. *IEEE Software*, 35(1), 82–88.
- [18] Kushvanth Chowdary Nagabhyru. (2023). Accelerating Digital Transformation with AI Driven Data Engineering: Industry Case Studies from Cloud and IoT Domains. *Educational Administration: Theory and Practice*, 29(4), 5898–5910. <https://doi.org/10.53555/kuvey.v29i4.10932>
- [19] Bottou, L. (2010). Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT 2010* (pp. 177–186). Physica-Verlag.
- [20] Deep Learning-Driven Optimization of ISO 20022 Protocol Stacks for Secure Cross-Border Messaging. (2024). *MSW Management Journal*, 34(2), 1545-1554.
- [21] Burns, B., Beda, J., & Hightower, K. (2019). *Kubernetes: Up & running* (2nd ed.). O'Reilly Media.
- [22] Cao, Y., Jia, X., Chen, Y., Lin, S., & Zhang, X. (2020). Deep learning for industrial inspection: A survey. *IEEE Transactions on Industrial Informatics*, 16(8), 4876–4891.
- [23] Meda, R. (2023). Intelligent Infrastructure for Real-Time Inventory and Logistics in Retail Supply Chains. *Educational Administration: Theory and Practice*.
- [24] Chen, D., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794). ACM.
- [25] Chen, M., Mao, S., & Liu, Y. (2014). Big data: A survey. *Mobile Networks and Applications*, 19(2), 171–209.
- [26] Aitha, A. R. (2023). CloudBased Micro services Architecture for Seamless Insurance Policy Administration. *International Journal of Finance (IJFIN)-ABDC Journal Quality List*, 36(6), 607-632.
- [27] Codd, E. F. (1970). A relational model of data for large shared data banks. *Communications of the ACM*, 13(6), 377–387.
- [28] Collins, E., & Nechvatal, J. (2020). NIST privacy framework: A tool for improving privacy through enterprise risk management (Version 1.0). National Institute of Standards and Technology.
- [29] Segireddy, A. R. (2024). Machine Learning-Driven Anomaly Detection in CI/CD Pipelines for Financial Applications. *Journal of Computational Analysis and Applications*, 33(8).
- [30] Craswell, N., Mitra, B., Yilmaz, E., Campos, D., & Voorhees, E. M. (2020). Overview of the TREC 2020 Deep Learning Track. In *Proceedings of the Text REtrieval Conference (TREC 2020)*. NIST.
- [31] Croft, W. B., Metzler, D., & Strohman, T. (2010). *Search engines: Information retrieval in practice*. Addison-Wesley.
- [32] Dai, Z., Yang, Z., Yang, Y., Cohen, W. W., Carbonell, J., Le, Q. V., & Salakhutdinov, R. (2019). Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (pp. 2978–2988). ACL.
- [33] Varri, D. B. S. (2024). Adaptive and Autonomous Security Frameworks Using Generative AI for Cloud Ecosystems. Available at SSRN 5774785.
- [34] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT 2019* (pp. 4171–4186). ACL.
- [35] Ding, S. X. (2014). *Data-driven design of fault diagnosis and fault-tolerant control systems*. Springer.
- [36] Singireddy, J. (2024). AI-Enhanced Tax Preparation and Filing: Automating Complex Regulatory Compliance.

- European Data Science Journal (EDSJ) p-ISSN 3050-9572 en e-ISSN 3050-9580, 2(1).
- [37] Dourish, P. (2004). What we talk about when we talk about context. *Personal and Ubiquitous Computing*, 8(1), 19–30.
- [38] Dwork, C., & Roth, A. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4), 211–407.
- [39] Keerthi Amistapuram. (2024). Federated Learning for Cross-Carrier Insurance Fraud Detection: Secure Multi-Institutional Collaboration. *Journal of Computational Analysis and Applications (JoCAAA)*, 33(08), 6727–6738. Retrieved from <https://www.eudoxuspress.com/index.php/pub/article/view/3934>
- [40] Evans, D. (2011). *The Internet of Things: How the next evolution of the internet is changing everything*. Cisco Internet Business Solutions Group.
- [41] Farooq, M. S., Khan, Z., Ahmad, R., Islam, S. U., & Kim, S. W. (2023). A survey on the role of industrial IoT in manufacturing for Industry 4.0. *Sensors*, 23(21), 8958.
- [42] Varri, D. B. S. (2023). *Advanced Threat Intelligence Modeling for Proactive Cyber Defense Systems*. Available at SSRN 5774926.
- [43] Fowler, M. (2018). *Refactoring: Improving the design of existing code* (2nd ed.). Addison-Wesley.
- [44] Gandomi, A., & Haider, M. (2015). Beyond the hype: Big data concepts, methods, and analytics. *International Journal of Information Management*, 35(2), 137–144.
- [45] Paleti, S. (2024). Transforming Financial Risk Management with AI and Data Engineering in the Modern Banking Sector. *American Journal of Analytics and Artificial Intelligence (ajaaai)* with ISSN 3067-283X, 2(1).
- [46] Grieves, M., & Vickers, J. (2017). Digital twin: Mitigating unpredictable, undesirable emergent behavior in complex systems. In F.-J. Kahlen, S. Flumerfelt, & A. Alves (Eds.), *Transdisciplinary perspectives on complex systems* (pp. 85–113). Springer.
- [47] Sheelam, G. K., & Koppolu, H. K. R. (2024). From Transistors to Intelligence: Semiconductor Architectures Empowering Agentic AI in 5G and Beyond. *Journal of Computational Analysis and Applications (JoCAAA)*, 33(08), 4518-4537.
- [48] Gray, J., & Reuter, A. (1993). *Transaction processing: Concepts and techniques*. Morgan Kaufmann.
- [49] Garapati, R. S. (2023). *Optimizing Energy Consumption in Smart Build-ings Through Web-Integrated AI and Cloud-Driven Control Systems*.
- [50] Guo, J., Fan, Y., Ai, Q., & Croft, W. B. (2020). A deep look into neural ranking models for information retrieval. *Information Processing & Management*, 57(6), 102067.
- [51] Han, S., Mao, H., & Dally, W. J. (2016). Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. In *Proceedings of ICLR 2016*.
- [52] Inala, R. *Revolutionizing Customer Master Data in Insurance Technology Platforms: An AI and MDM Architecture Perspective*.
- [53] He, W., Xu, L. D., & Chen, H. (2014). Internet of Things in industries: A survey. *IEEE Transactions on Industrial Informatics*, 10(4), 2233–2243.
- [54] Varri, D. B. S. (2022). A Framework for Cloud-Integrated Database Hardening in Hybrid AWS-Azure Environments: Security Posture Automation Through Wiz-Driven Insights. *International Journal of Scientific Research and Modern Technology*, 1(12), 216-226.
- [55] Hohpe, G., & Woolf, B. (2003). *Enterprise integration patterns: Designing, building, and deploying messaging solutions*. Addison-Wesley.
- [56] Amistapuram, K. (2024). Generative AI in Insurance: Automating Claims Documentation and Customer Communication. *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*, 15(3), 461–475. <https://doi.org/10.61841/turcomat.v15i3.15474>
- [57] IEC. (2018). IEC 62443-3-3:2013 + AMD1:2017 + AMD2:2020 Industrial communication networks—Network and system security—Part 3-3: System security requirements and security levels. International Electrotechnical Commission.
- [58] ISO. (2018). ISO/IEC 27001:2018 Information security management systems—Requirements. International Organization for Standardization.
- [59] Guntupalli, R. (2024). Enhancing Cloud Security with AI: A Deep Learning Approach to Identify and Prevent Cyberattacks in Multi-Tenant Environments. Available at SSRN 5329132.
- [60] IT Governance Institute. (2012). *COBIT 5: A business framework for the governance and management of enterprise IT*. ISACA.
- [61] Järvelin, K., & Kekäläinen, J. (2002). Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems*, 20(4), 422–446.
- [62] Koppolu, H. K. R., & Sheelam, G. K. (2024). Machine Learning-Driven Optimization in 6G Telecommunications: The Role of Intelligent Wireless and Semiconductor Innovation. *Global Research Development (GRD)* ISSN: 2455-5703, 9(12).
- [63] Johnson, J., Douze, M., & Jégou, H. (2019). Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*,

7(3), 535–547.

- [64] Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., D'Oliveira, R. G. L., Rouayheb, S. E., Gascón, A., Ghazi, B., Gibbons, P. B., Hastie, T., Hazy, T., Kalenichenko, D., Kamath, G., ... Zhao, S. (2021). Advances and open problems in federated learning. *Foundations and Trends in Machine Learning*, 14(1–2), 1–210.
- [65] Lahari Pandiri, "AI-Powered Fraud Detection Systems in Professional and Contractors Insurance Claims," *International Journal of Innovative Research in Electrical, Electronics, Instrumentation and Control Engineering (IJREEICE)*, DOI 10.17148/IJREEICE.2024.121206.
- [66] Katz, R., Goldschmidt, T., & Grady, J. (2021). Edge computing security: A survey. *IEEE Access*, 9, 158820–158840.
- [67] Khattab, O., & Zaharia, M. (2020). ColBERT: Efficient and effective passage search via contextualized late interaction over BERT. In *Proceedings of SIGIR 2020* (pp. 39–48). ACM.
- [68] Rongali, S. K. (2023). Explainable Artificial Intelligence (XAI) Framework for Transparent Clinical Decision Support Systems. *International Journal of Medical Toxicology and Legal Medicine*, 26(3), 22-31.
- [69] Lee, J., Bagheri, B., & Kao, H.-A. (2015). A cyber-physical systems architecture for Industry 4.0-based manufacturing systems. *Manufacturing Letters*, 3, 18–23.
- [70] Lee, J., Jin, C., & Bagheri, B. (2017). Cyber physical systems for predictive production systems. *Production Engineering*, 11(2), 155–165.
- [71] Inala, R. AI-Powered Investment Decision Support Systems: Building Smart Data Products with Embedded Governance Controls.
- [72] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-T., Rocktäschel, T., Riedel, S., & Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33, 9459–9474.
- [73] Mashetty, S., Challa, S. R., ADUSUPALLI, B., Singireddy, J., & Paleti, S. (2024). Intelligent Technologies for Modern Financial Ecosystems: Transforming Housing Finance, Risk Management, and Advisory Services Through Advanced Analytics and Secure Cloud Solutions. *Risk Management, and Advisory Services Through Advanced Analytics and Secure Cloud Solutions* (December 12, 2024).
- [74] Liu, T.-Y. (2009). Learning to rank for information retrieval. *Foundations and Trends in Information Retrieval*, 3(3), 225–331.
- [75] Rongali, S. K., & Kumar Kakarala, M. R. (2024). Existing challenges in ethical AI: Addressing algorithmic bias, transparency, accountability and regulatory compliance.
- [76] Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765–4774.
- [77] Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to information retrieval*. Cambridge University Press.
- [78] Guntupalli, R. (2024). AI-Powered Infrastructure Management in Cloud Computing: Automating Security Compliance and Performance Monitoring. Available at SSRN 5329147.
- [79] Mell, P., & Grance, T. (2011). The NIST definition of cloud computing (NIST SP 800-145). National Institute of Standards and Technology.
- [80] Nagubandi, A. R. (2023). Advanced Multi-Agent AI Systems for Autonomous Reconciliation Across Enterprise Multi-Counterparty Derivatives, Collateral, and Accounting Platforms. *International Journal of Finance (IJFIN)-ABDC Journal Quality List*, 36(6), 653-674.
- [81] Mohri, M., Rostamizadeh, A., & Talwalkar, A. (2018). *Foundations of machine learning* (2nd ed.). MIT Press.
- [82] Keerthi Amistapuram. (2023). Privacy-Preserving Machine Learning Models for Sensitive Customer Data in Insurance Systems. *Educational Administration: Theory and Practice*, 29(4), 5950–5958. <https://doi.org/10.53555/kuvey.v29i4.10965>
- [83] NIST. (2020). Security and privacy controls for information systems and organizations (NIST SP 800-53 Rev. 5). U.S. Department of Commerce.
- [84] Chava, K. (2024). The Role of Cloud Computing in Accelerating AI-Driven Innovations in Healthcare Systems. *European Advanced Journal for Emerging Technologies (EAJET)-p-ISSN 3050-9734 en e-ISSN 3050-9742*, 2(1).
- [85] Object Management Group. (2016). Business process model and notation (BPMN), version 2.0.2. OMG.
- [86] Object Management Group. (2019). Decision model and notation (DMN), version 1.3. OMG.
- [87] Rongali, S. K. (2024). Federated and Generative AI Models for Secure, Cross-Institutional Healthcare Data Interoperability. *Journal of Neonatal Surgery*, 13(1), 1683-1694.
- [88] Pan, Y., Zhang, L., & Liu, S. (2022). Data-driven quality prediction and anomaly detection in smart manufacturing: A review. *Journal of Manufacturing Systems*, 63, 53–72.
- [89] AI and ML-Driven Optimization of Telecom Routers for Secure and Scalable Broadband Networks. (2024). *MSW Management Journal*, 34(2), 1145-1160.
- [90] Singh, R., Auluck, N., & Rana, O. (2023). Edge AI: A survey. *Results in Engineering*, 18, 101053.